

Running head: MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Moral values in causal attribution:  
Evidence from the implicit verb causality task and explicit judgments**

6429 words without references, abstract and notes

**Moral values in causal attribution:**

**Evidence from the implicit verb causality task and explicit judgments**

Laura Niemi\*, University of Toronto

Joshua Hartshorne, Boston College

Tobias Gerstenberg, Stanford University

Matthew Stanley, Duke University

Liane Young, Boston College

6429 words without references, abstract and notes

**\*Corresponding author:**

Laura Niemi

University of Toronto

Munk School of Global Affairs and Public Policy

315 Bloor Street West

Toronto, ON M5S 1W7

[laura.niemi@utoronto.ca](mailto:laura.niemi@utoronto.ca)

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Abstract**

Prior work found that “binding values,” moral values that protect the group, are linked to victim-blaming. It was speculated that binding values involve an understanding of the causal structure of harmful events in which causation is placed on harmed people. The present research used the implicit causality task from psycholinguistics (e.g., “Bob *verb*-ed Amy because... he or she?”), explicit judgments of causal contributions, and measurement of participants’ moral values to investigate how moral values relate to interpretation of the causal structure of events. Using two verb sets and two independent replications ( $N = 459$ ,  $N = 249$ ,  $N = 788$ ), we found that binding values predicted selection of the object (victim) as the cause in the implicit causality task for harmful events. Binding values also predicted explicit causal attributions. The findings indicate that moral values that support close social bonds reliably predict causal attributions to people affected by harmful events.

*Keywords:* Morality; Social Cognition; Cognition; Causality; Psycholinguistics

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Moral values in causal attribution:****Evidence from the implicit verb causality task and explicit judgments**

It is no secret that people often disagree about what's right or wrong. But when a negative event happens, people with very different moral values (Graham et al., 2011; Niemi & Young, 2016) often ask the same sorts of questions, such as: *How did this happen?*, *Who could have let it happen?*, and *Would this have happened to somebody else, regardless?* (Alicke, 2000; Alicke, Mandel, Hilton, Gerstenberg & Lagnado, 2015; Heider, 1958; Malle, Guglielmo & Monroe, 2014). These questions reveal that, disagreements aside, everyone engages causal cognition when assigning blame and responsibility, as indicated by much work in moral psychology (Alicke, 2000; Alicke et al, 2015; Cushman, 2008; Heider, 1958; Malle et al, 2014; Shaver, 1985). At what point in the evaluation of events does moral disagreement emerge?

One explanation is that moral disagreement is largely a function of differences of opinion about the specific actions that should be considered moral or immoral (e.g., premarital sex, abortion). Research has shown that there are, in fact, clear patterns in what people do or do not find immoral. In particular, there are at least two clusters of moral values (Graham et al., 2011): (i) caring and fairness values, deemed "individualizing values" because they are extended to each individual regardless of group membership; and (ii) loyalty, obedience, and purity values, deemed "binding values" because they purportedly keep people bound into relationships and groups.

These two clusters of values could lead to moral disagreement not only because of differences in what is regarded as immoral, but because they involve very different perspectives on cause and effect. Violations in the first cluster (individualizing values) involve an agent causing harm to a patient (the "moral dyad"— Gray, Young & Waytz, 2012). Violations in the second cluster may not involve any obvious harm done by an agent to a patient. Indeed, they don't require clear boundaries between the roles of agent and patient, may require a third person, and can sometimes be performed solely by an agent (Niemi & Young, 2016). Moreover,

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 the same event (*A killed B*) may be immoral from the perspective of individualizing values but  
4  
5 morally *obligatory* from the perspective of binding values (*A was ordered to kill the traitor B*).  
6  
7 This shift in moral acceptability is made possible by viewing the killing from an alternative causal  
8  
9 perspective in which *A*'s causal contribution as the agent is reduced and *B*'s contribution as the  
10  
11 patient is increased.  
12

13  
14 People tend to systematically differ in their endorsement of binding and individualizing  
15  
16 values (Graham et al., 2011): binding values are higher in people who more strongly endorse  
17  
18 conservative political ideology. Previously, Niemi and Young (2016) found increased victim-  
19  
20 blaming in people higher in binding values even after controlling for political orientation. Ratings  
21  
22 of victims' causal responsibility mediated the relationship between victim-blaming and binding  
23  
24 values, in line with work showing that causal judgments feed judgments of blame (e.g., Malle et  
25  
26 al., 2012). It was speculated that binding values might involve a different understanding of the  
27  
28 causal structure of events of harm and force, leading people higher in these values to be more  
29  
30 likely to shift causation off of harm-doers and over to harmed people. The present research  
31  
32 directly investigates this possibility.  
33

34  
35 We first investigate this hypothesis using a task that measures intuitions about likely  
36  
37 causes for events from psycholinguistics, the implicit causality (IC) task (Brown & Fish, 1983;  
38  
39 Garvey & Caramazza, 1974; Hartshorne & Snedeker, 2012; Rudolph & Forsterling, 1997). In  
40  
41 this task, participants read sentences such as  
42

43 (a) Bob murdered Amy because...

44  
45 and chose whether to continue the sentence with a pronoun referring to the agent ("he")  
46  
47 or patient ("she"). Continuing (a) with *he* or *she* reveals an expectation that the murder is due to  
48  
49 something that the subject (Bob) did or the object (Amy) did, respectively. Implicit causality  
50  
51 research indicates that this choice tends to vary by verb (*praise*, *frighten*, etc.), with some verbs  
52  
53 reliably prompting selection of pronouns referring to the subject (e.g., "subject biased" verbs like  
54  
55 *frighten*) and some prompting selection of pronouns referring to the object (e.g., "object biased"  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 verbs like *praise*), suggesting that people have systematic expectations about how some  
4  
5 categories of events came about (Brown & Fish, 1983; Garvey & Caramazza, 1974; Hartshorne  
6  
7 & Snedeker, 2012; Rudolph & Forsterling, 1997; Bott & Solstad, 2014; Ferstl, Garnham &  
8  
9 Manouilidou, 2011; Hartshorne, 2013; Pickering & Majid, 2007; Rudolph, 2008). Indeed, the only  
10  
11 well-established predictor of IC is verb semantics. Specifically, when verbs have been clustered  
12  
13 into classes based on fine-grained analyses of shared semantic and syntactic features — which  
14  
15 include information about causation — they tend to share implicit causality biases to the subject  
16  
17 or object (Hartshorne, 2013; Hartshorne & Snedeker, 2012; Kipper-Schuler, 2006).

18  
19  
20 Interestingly, although many researchers have suggested that IC is broadly affected by  
21  
22 an individual's beliefs about the world, such as perceived social hierarchy or gender roles, the  
23  
24 accompanying evidence has been inconsistent (Garvey & Caramazza, 1974; Bott & Solstad,  
25  
26 2014; Ferstl, Garnham & Manouilidou, 2011; Hartshorne, 2013; Pickering & Majid, 2007).  
27  
28 Despite the potential for IC to be affected by relevant individual differences (e.g., the role of  
29  
30 moral values on morally relevant verbs' IC), this topic is largely unexplored – with the exception  
31  
32 of work examining gender differences in pronoun comprehension and interpretation (Arnold,  
33  
34 2015).

35  
36  
37 Since explanations typically point to the most relevant or contributory causes (e.g.,  
38  
39 Hilton, 1990; Hesslow, 1988; Lombrozo, 2006 ), IC selections provide a window into how people  
40  
41 think about the causes of events. For instance, note that explaining (a) in terms of Bob is more  
42  
43 consistent with the moral dyad framework, in which active perpetrators harm passive victims.  
44  
45 The present research leveraged the IC task to investigate the causal structure underlying the  
46  
47 blame and responsibility judgments of people ranging in moral values, and also let us contribute  
48  
49 to the understanding of individual differences in IC responses. We investigated whether people  
50  
51 high in binding values were more likely to exhibit an object-bias in the IC task, expecting  
52  
53 explanations of harm events to focus on the *victim*, given prior, repeatedly replicated findings  
54  
55 that binding values predict victim blame and stigmatization (Niemi & Young, 2016). Less  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 consistent prior evidence linking individualizing values and perpetrator blame suggested that  
4 individualizing values would be linked to the opposing pattern, subject-bias for harm events. We  
5 also tested responses to morally irrelevant, neutral verbs to rule out the possibility that people  
6 high in binding values were generally more likely to consider affected people to be causal  
7 contributors.  
8  
9  
10  
11  
12

13  
14 We also measured participants' explicit causal judgments, including judgments about  
15 whether the agent's action was necessary and sufficient for the outcome, and whether the  
16 patient allowed, controlled, and deserved what happened. Because we expected individuals  
17 who prioritize binding values to be less likely to apply the moral dyad framework in which agents  
18 are causal and blameworthy ("agent-harmed-patient") when reasoning about immoral events,  
19 we expected these participants to view agents as less necessary or sufficient and to view  
20 patients as more likely to have allowed, controlled or deserved the events. Moreover, we  
21 expected that IC object-bias would be directly related to judgments that agents were less  
22 necessary and sufficient and that patients allowed, controlled, and deserved the events.  
23  
24  
25  
26  
27  
28  
29  
30  
31

32  
33 We measured participants' sensitivity to victim suffering (how "injured" participants  
34 considered victims) and stigmatization of victims (how "contaminated" participants considered  
35 victims) to understand how these explicit morally-relevant attitudes about harmed people  
36 (sensitivity vs. stigmatization) related to implicit causality selections, and to test replication of  
37 prior work. In prior work (Niemi & Young, 2016), increased sensitivity to victim suffering – rating  
38 victims as more "injured" – was associated with higher individualizing values. Increased  
39 stigmatization of victims – rating victims as more "contaminated", like victim-blaming, was  
40 associated with higher binding values. Victim injury ratings were also negatively correlated with  
41 victim contamination ratings, suggesting that viewing victims as contaminated is inconsistent  
42 with viewing them as passive victims of harm. Finally, we measured moral values with the Moral  
43 Foundations Questionnaire, which has been used extensively in prior work to measure people's  
44 diverse moral values (e.g., Graham et al., 2011; Niemi & Young, 2016).  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Roadmap**

Here, we outline the four tested hypotheses in the present research. The first three hypotheses examine how moral values are related to responses on the implicit causality task, explicit causal judgments, and the propensity to stigmatize victims (or be sensitive to victim suffering). The last hypothesis examines their interrelationships. These are tested in the primary Study 1, and two replication datasets, as indicated below.

- (1) We expected those higher in binding values to be more likely to select the object over the subject as the referent (“object-bias”) for harm and force events, but not for neutral events in the IC task.
- (2) We expected binding values to be negatively related to participants’ explicit causal judgments of the agent’s necessity and sufficiency for harm and force events, and positively related to judgments of the patient’s capacity to allow, control and deserve harm and force events.
- (3) We expected binding values to positively correlate with stigmatization of victims, and individualizing values to positively correlate with sensitivity to victim suffering.
- (4) We hypothesized that IC object-bias for harm and force events would be related to reduced sensitivity to victim suffering and judgments of agents as less necessary and sufficient, greater stigmatization of victims and increased judgments that patients allowed, controlled and deserved harm and force events.

Study 1 tests all four of these hypotheses. We attempt to replicate all findings involving the IC bias from Study 1 in Replication Dataset 1. We attempt to replicate all findings involving the IC bias again using an expanded set of verbs and a larger sample size in Replication Dataset 2.



## MORAL VALUES AND CAUSAL ATTRIBUTION

## Study 1

## Method

Participants ( $N = 459$ ) were recruited online via Amazon's Mechanical Turk ( $M_{age} = 37.25$  years,  $SD_{age} = 31.39$ ; 207 selected female, 247 selected male, 5 selected other or missing). We excluded 189 additional individuals who failed attention checks.<sup>1</sup> We aimed to have approximately 200 participants in each condition (*Male-verb-ed-female versus Female-verb-ed-male*, described below) in line with past work showing that associations among moral values, blame, and responsibility were found in samples of approximately this size (Niemi & Young, 2016). We also collected data from additional participants online<sup>2</sup> via Amazon's Mechanical Turk and attempted to replicate effects obtained in the primary experiment with Replication Dataset 1 ( $N = 249$ ) ( $M_{age} = 35.87$  years,  $SD_{age} = 13.49$ ; 114 selected female, 133 selected male, 2 selected other or missing) and Replication Dataset 2 ( $N = 788$ ) ( $M_{age} = 36.32$  years,  $SD_{age} = 12.88$ ; 279 selected female, 504 selected male, 5 selected other or missing). Data and materials are available at <https://github.com/BLINDEDFORREVIEW/>. Methodological differences between Study 1 and the Replication Datasets are described in the **Supplementary Materials**. The institutional review board at Boston College approved all studies, and informed consent was obtained via an online form from all participants.

Moral values in the five foundations (caring, fairness, loyalty, obedience to authority, and purity) were assessed using the 30-item Moral Foundations Questionnaire (MFQ; Graham et al, 2011). Participants also provided demographic information including political orientation, gender, and religiosity, and they completed the Ambivalent Sexism Inventory (Glick & Fiske, 1996; the present analyses do not involve the ASI).

The implicit causality task involved 24 minimal event descriptions in the form: "[*Subject*] *verb-ed* [*Object*] *because...*" – e.g., "*Bob coerced Amy because...*" – with half of the participants receiving male sentence subjects and female sentence objects, and vice versa for the other half

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 in order to equalize gender of the person in the subject and object positions. Participants were  
4 asked to “Please select which word you think would follow.” They were offered the choices “*he*”  
5 or “*she*” (counterbalanced order across items). Verbs described highly morally-relevant events  
6 (the “harm/force” verbs, henceforth), and neutral events (“neutral filler” verbs; see **Table 1** for  
7 verbs).<sup>3</sup> Note that we purposefully selected verbs describing events likely to be of importance for  
8 informing theories of morality (e.g., *kill* and *rape*). We conducted analyses in which we  
9 examined the effect of harm/force verbs and neutral filler verbs using linear mixed-effects  
10 models.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Table 1.** Verbs used in analyses in Study 1, Replication 1 and 2.

<b>Study 1 and Replication 1</b>		<b>Replication 2</b>			
<b><i>Harm/Force</i></b>		<b><i>Harm/Force</i></b>			
clobbered	manipulated	assaulted	forced	molested	silenced
coerced	raped	clobbered	groped	persuaded	spanked
enslaved	robbed	coerced	impaled	pressured	stabbed
forced	stabbed	convinced	influenced	raped	strangled
influenced	strangled	enslaved	killed	robbed	tempted
killed	tempted	enticed	manipulated	seduced	walloped
<b><i>Neutral Filler</i></b>		<b><i>Neutral Filler</i></b>			
approached	praised	appraised	complimented	honored	quoted
congratulate	quoted	approached	congratulated	impressed	raced
delighted	skipped	boggled	delighted	massaged	skipped
impressed	thanked	caressed	diverted	observed	thanked
observed	transported	celebrated	fondled	praised	tickled
		comforted	greeted	puzzled	transported

## MORAL VALUES AND CAUSAL ATTRIBUTION

To assess explicit causal judgments, we gathered participants' beliefs about agents' and patients' causal contributions. After completing all the implicit causality task items, participants viewed the same events they had seen in the implicit causality block without the "because" connective (e.g., "*Bob coerced Amy.*"). They were asked to "weigh the following possibilities" in the following order:

1. Agent Unnecessary: e.g., "Would [*Amy*] have been [*coerced*] by someone else?"
2. Agent Sufficient: e.g., "Would [*Bob*] [*coerce*] someone else?"
3. Patient Control: e.g., "Did [*Amy*] have control over the occurrence of the event?"
4. Patient Allowing: e.g., "Did [*Amy*] let the event happen?"
5. Patient Desert: e.g., "Could [*Amy*] have deserved the event?"

Participants responded using sliding scales (0 = "*Definitely No*", 50 = "*Unsure*", 100 = "*Definitely Yes*").

Finally, we measured sensitivity versus stigmatization toward victims as in prior work (Niemi & Young, 2016). We asked participants to rate in counterbalanced order how "*contaminated / tainted*" and "*injured / wounded*" they considered hypothetical crime victims (crimes: *molestation, rape, strangling, stabbing*) on a sliding scale from 0 (*Not at all*) to 7 (*Very much*). As in prior work, only these four crimes were used to obtain measures of how "*contaminated / tainted*" and "*injured / wounded*" participants rated hypothetical crime victims. Average ratings of victims across events as contaminated/tainted and injured/wounded served as indices of stigmatization of victims and sensitivity to victim suffering, respectively.

*Statistical Analyses.* Data were analyzed in several ways to address different questions. First, using R (R Development Core Team, 2009) with the lme4 software package (Bates, Maechler, Bolker, & Walker, 2015), we computed a series of generalized linear mixed-effects

## MORAL VALUES AND CAUSAL ATTRIBUTION

models. For models with binary outcome variables, significance and 95% CIs around beta-estimates were computed using Wald tests. For models with non-binary outcome variables, significance for fixed effects was assessed using Satterthwaite approximations to degrees of freedom, and 95% CIs around beta-estimates were computed using parametric bootstrapping. In all models, participant and verb were included as crossed random effects (random intercepts only).<sup>4</sup> Finally, to address questions about the relationship between moral values, stigmatization, and sensitivity to victim suffering, we computed a series of Spearman's rank-order correlations.

**Results and Discussion**

We address our four hypotheses in the order that they were presented in the Introduction.

*Moral Values and Implicit Causality Object-bias.* We expected those higher in binding values to be more likely to select the object over the subject as the referent ("object-bias") for harm and force events, but not for neutral events. We tested the relationship between moral values and implicit causality object-bias with a series of generalized linear mixed-effects models (link = "logit"). First, a generalized linear mixed-effects regression model was computed in which verb type (harm/force (coded as 0) *versus* neutral filler (coded as 1)) and binding values were included as fixed predictors of the propensity to select the object (coded as 1) *relative to* the subject (coded as 0) as the referent. There was a significant interaction between verb type and binding values in Study 1 and Replication Datasets (see **Table 2**). To further interrogate these significant interaction effects, generalized linear mixed-effects models were computed for harm/force verbs and neutral filler verbs, taken separately.

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Table 2.** *The results of two generalized linear mixed-effects regression models—each with verb type and binding values as predictors of selecting object relative to the subject as the referent.*

	<i>b</i>	<i>SE</i>	<i>Z</i>	<i>p</i>	95% CI
<b>Study 1</b>					
Verb Type	2.17	.43	5.06	< .0001	[1.33, 3.02]
Binding Values	.37	.05	7.02	< .0001	[.27, .48]
Verb Type x Binding Values	-.27	.05	-5.60	< .0001	[-.36, -.18]
<b>Replication Dataset 1</b>					
Verb Type	2.43	.47	5.17	< .0001	[1.51, 3.36]
Binding Values	.43	.08	5.27	< .0001	[.27, .59]
Verb Type x Binding Values	-.35	.07	-4.76	< .0001	[-.49, -.20]
<b>Replication Dataset 2</b>					
Verb Type	1.28	.32	4.35	< .0001	[.75, 2.00]
Binding Values	.19	.04	4.66	< .0001	[.11, .27]
Verb Type x Binding Values	-.10	.03	-2.54	.0004	[-.15, -.04]

*Note.* Study 1 ( $N = 459$ ), Replication Dataset 1 ( $N = 249$ ), Replication Dataset 2 ( $N = 788$ ). All 95% CIs are for the beta-estimates.

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3  
4  
5  
6  
7  
8 For harm/force verbs only, a generalized linear mixed-effects regression model was  
9 computed for which binding values was included as the fixed predictor of the propensity to  
10 select the object (coded as 1) *relative to* the subject (coded as 0) as the referent. This analysis  
11 yielded a significant effect of binding values on the likelihood of selecting the object as the  
12 referent ( $p < .0001$ ). We obtained the same pattern of results in Replication Dataset 1 ( $p <$   
13  $.0001$ ) and Replication Dataset 2 ( $p < .0001$ ). In all three datasets, participants higher in binding  
14 values were more likely to select the object over the subject as the referent (object-bias) for  
15 harm and force events.  
16  
17  
18  
19  
20  
21  
22  
23

24 Importantly, for the neutral filler verbs, there was no significant effect of binding values  
25 on selection of the object *relative to* the subject as the referent in the primary study or  
26 Replication Dataset 1 (both  $ps > .05$ ). For the neutral filler verbs in Replication Dataset 2, there  
27 was a small but significant effect of binding values on the selection of the object *relative to* the  
28 subject as the referent ( $p = .02$ ). So, the effect was much larger for harm/force verbs than for  
29 neutral filler verbs.  
30  
31  
32  
33  
34  
35  
36

37 Next, we tested a number of additional considerations related to the implicit causality  
38 object-bias. All of these findings are presented in full in **Supplementary Materials**. First, given  
39 that prior work has identified relationships between binding values and political orientation,  
40 gender, and religiosity (Graham et al., 2011), we wanted to ensure that binding values predicted  
41 the implicit causality object-bias above and beyond these other variables. Analyses showed that  
42 binding values remain consistent significant predictors of implicit causality object-bias for  
43 harm/force verbs after controlling for political orientation, gender, and religiosity in the three  
44 datasets. Second, we found no effect of individualizing values on the propensity to select the  
45 object *relative to* the subject as the referent for harm/force verbs or neutral filler verbs. Third,  
46 gender condition (male-verbed-female *versus* female-verbed-male) was related to the implicit  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

causality object-bias for harm/force verbs. Specifically, in all datasets, participants were more likely to select men for harm/force events. However, binding values continued to significantly predict the implicit causality object-bias despite the gender effect.

Overall, the implicit causality results support our first hypothesis. Participants higher in binding values were more likely to select the object over the subject as the referent for harm and force events, but not for neutral events. These effects held after controlling for a variety of other variables.

*Explicit Causal Judgments: Agents' and Patients' Contributions.* We next tested our hypothesis that binding values would be negatively related to participants' judgments about agents' necessity and sufficiency, and positively related to their judgments of patients' capacities to allow, control, and deserve events of harm and force; as a comparison, we investigated whether the opposite patterns would be observed in the case of individualizing values. We first computed a series of linear mixed-effects models in which binding values and verb type (harm/force (coded as 0) versus neutral filler (coded as 1)) were included as fixed predictors of judgments for necessity, sufficiency, allowing, controlling, and deserving (in separate models). In all five models, there was a significant interaction effect between binding values and verb type (see **Table 3**). To further interrogate these significant interaction effects, linear mixed-effects models were computed for harm/force verbs and neutral filler verbs, taken separately.



## MORAL VALUES AND CAUSAL ATTRIBUTION

**Table 3.** The results of five different linear mixed-effects regression models are depicted. In all models, verb type and binding values were fixed predictors; necessity, sufficiency, allowing, controlling, and deserving were the outcome variables in the different models.

	<i>b</i>	<i>SE</i>	<i>t</i>	<i>p</i>	95% CI
<b>Outcome: Necessity</b>					
Verb Type	-33.83	4.70	-7.20	< .0001	[-42.77, -25.96]
Binding Values	-3.36	.78	-4.33	< .0001	[-4.75, -1.86]
Verb Type x Binding Values	2.96	.40	7.37	< .0001	[2.21, 3.74]
<b>Outcome: Sufficiency</b>					
Verb Type	-3.32	1.87	-1.70	.095	[-7.13, .62]
Binding Values	-1.77	.83	-2.14	.033	[-3.55, -.16]
Verb Type x Binding Values	1.18	.35	3.40	.0007	[.44, 1.87]
<b>Outcome: Allow</b>					
Verb Type	30.95	6.51	4.76	< .0001	[18.59, 44.02]
Binding Values	5.71	.77	7.42	< .0001	[4.16, 7.37]
Verb Type x Binding Values	-1.71	.49	-3.52	< .0001	[-2.61, -.71]
<b>Outcome: Control</b>					
Verb Type	18.30	5.90	3.10	.005	[7.89, 30.81]
Binding Values	4.08	.74	5.54	< .0001	[2.70, 5.54]
Verb Type x Binding Values	-1.22	.48	-2.53	.012	[-2.22, -.39]
<b>Outcome: Deserve</b>					
Verb Type	47.38	5.58	8.49	< .0001	[35.69, 57.98]
Binding Values	3.83	.66	5.80	< .0001	[2.39, 5.11]
Verb Type x Binding Values	-3.23	.45	-7.17	< .0001	[-4.04, -2.34]

Note. All 95% CIs are for the beta-estimates.

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 For harm/force verbs only, five linear mixed-effects models with binding values as the  
4 fixed predictor of necessity, sufficiency, allowing, controlling, or deserving judgments (in  
5 separate models) revealed that binding values were negatively related to participants'  
6 judgments about the agent's necessity ( $p = .0002$ ) and sufficiency ( $p = .046$ ), and positively  
7 related to their judgments about the patient's capacity to allow ( $p < .0001$ ), control ( $p < .0001$ ),  
8 and deserve ( $p < .0001$ ) the events. Importantly, for the neutral filler verbs, there was no  
9 significant effect of binding values on judgments of necessity, sufficiency, or desert (all  $ps >$   
10  $.05$ ). However, there were significant effects of binding values on judgments of allowing ( $p <$   
11  $.0001$ ) and controlling ( $p = .0006$ ). Nevertheless, for allowing and controlling judgments, the  
12 magnitude of the effect was larger for harm/force verbs than for neutral filler verbs.

13  
14 For the purposes of comparison, we next computed a series of linear mixed-effects  
15 models in which individualizing values and verb type (harm/force (coded as 0) *versus* neutral  
16 filler (coded as 1)) were included as fixed predictors of judgments for necessity, sufficiency,  
17 allowing, controlling, and deserving (in separate models). In models with necessity, allowing,  
18 controlling, and deserving, there were significant interaction effects between individualizing  
19 values and verb type (see **Table 4**). For sufficiency judgments, there was only a significant main  
20 effect of individualizing values. To further interrogate the four significant interaction effects,  
21 linear mixed-effects models were computed for harm/force verbs and neutral filler verbs, taken  
22 separately.

23  
24 For harm/force verbs only, four linear mixed-effects models with individualizing values as  
25 the fixed predictor of necessity, allowing, controlling, or deserving judgments (in separate  
26 models) revealed that individualizing values were not significantly related to participants'  
27 judgments about the agent's necessity ( $p = .88$ ), but individualizing values were negatively  
28 related to judgments of the patient's capacity to allow ( $p = .048$ ), control ( $p = .0006$ ), and  
29 deserve ( $p = .001$ ) the events. For the neutral filler verbs, there was only a significant effect of  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 individualizing values on judgments of deserving ( $p = .016$ ). There were no significant effects for  
4  
5 judgments of necessity, allowing, or controlling for the neutral filler verbs (all  $ps > .05$ ).  
6

7 Overall, these results support our second hypothesis. For harm/force verbs, binding  
8  
9 values were negatively related to participants' explicit causal judgments about the agent's  
10  
11 necessity and sufficiency, and positively related to judgments about the patient's capacity to  
12  
13 allow, control and deserve the events. In contrast, for harm/force verbs, an opposing pattern  
14  
15 was observed with individualizing values: they were negatively related to judgments about the  
16  
17 patient's capacity to allow, control, and deserve the events.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

**Table 4.** The results of five different linear mixed-effects regression models are depicted. In all models, verb type and individualizing values were fixed predictors; necessity, sufficiency, allowing, controlling, and deserving were the different outcome variables in the different models.

	<i>b</i>	<i>SE</i>	<i>t</i>	<i>p</i>	95% CI
<b>Outcome: Necessity</b>					
Verb Type	-12.54	5.26	-2.38	.022	[-23.06, -2.04]
Individualizing Values	.20	1.14	.18	.859	[-1.93, 2.39]
Verb Type x Individualizing Values	-2.15	.59	-3.66	.0003	[-3.30, -.96]
<b>Outcome: Sufficiency</b>					
Verb Type	4.35	2.84	1.53	.127	[-1.20, 9.91]
Individualizing Values	4.64	1.20	3.87	.0001	[2.56, 7.02]
Verb Type x Individualizing Values	-.69	.51	-1.35	.178	[-1.74, .29]
<b>Outcome: Allow</b>					
Verb Type	11.90	7.10	1.68	.103	[-3.00, 26.27]
Individualizing Values	-2.65	1.18	-2.25	.025	[-4.71, -.40]
Verb Type x Individualizing Values	2.66	.71	3.75	.0002	[1.42, 3.89]
<b>Outcome: Control</b>					
Verb Type	4.68	6.55	.71	.480	[-8.51, 18.00]
Individualizing Values	-4.14	1.09	-3.78	.0002	[-6.21, -2.09]
Verb Type x Individualizing Values	1.91	.71	2.70	.007	[.56, 3.19]
<b>Outcome: Deserve</b>					
Verb Type	3.68	6.17	.60	.555	[-9.50, 15.34]
Individualizing Values	-4.16	.98	-4.24	< .0001	[-6.11, -2.27]
Verb Type x Individualizing Values	6.65	.66	10.12	< .0001	[5.39, 7.89]

Note. All 95% CIs are for the beta-estimates.

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3           *Sensitivity versus Stigmatization toward Victims.* First, we computed a series of  
4 correlations to test replication of previously observed relationships between binding values and  
5 stigmatization of victims, and individualizing values and sensitivity to victim suffering, for a  
6 subset of harmful events (*rape, strangling, stabbing*)<sup>5</sup>. We replicated prior findings (Niemi &  
7 Young, 2016) of a positive relationship between binding values and ratings of victims as  
8 contaminated, and a positive relationship between individualizing values and ratings of victims  
9 as injured (see **Table 5**).

10  
11           Next, to address the subsequent hypotheses that stigmatization and sensitivity for  
12 victims might be related to implicit causality object-bias and explicit causal judgments (agents'  
13 and patients' contributions), we calculated object-bias for harm/force verbs and object-bias for  
14 neutral filler verbs by taking the probability of selecting the object as referent across the  
15 harm/force events and neutral filler events, respectively. Thus, "harm/force object-bias"  
16 represented each participant's tendency toward selecting the object over the subject (akin to an  
17 "implicit victim-blaming" score). Correspondingly, the "neutral filler object-bias" represented a  
18 tendency to select the object over the subject across events that do not involve harm and force.  
19 Additionally, we created an "*Agent Contribution*" aggregate variable by averaging the agent  
20 unnecessary ratings (reverse-coded) and agent sufficiency ratings, and a "*Patient Contribution*"  
21 aggregate variable by averaging patient control, patient allowing, and patient deserving ratings.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

41           A series of correlations indicated that, as hypothesized, ratings of victims as  
42 contaminated were significantly associated with a more pronounced implicit causality object-  
43 bias, increased patient contribution ratings, and decreased agent contribution ratings. By  
44 contrast, ratings of victims as injured were significantly negatively associated with implicit  
45 causality object-bias and with ratings of patients as causal contributors. They were also  
46 significantly positively associated with agent contribution ratings (**Table 5**). It is notable that  
47 people's ratings of how "contaminated" and "injured" they considered generic, unnamed victims  
48 of crimes (i.e., rape, stabbing, strangling) – completed in a separate part of the study – showed  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 reliable relationships with implicit causality object-bias — i.e., selection of the object as the  
4 cause for various harm and force events (such as: “Bob killed Amy because...*she*”).  
5  
6

7 Overall, these results largely supported our third and fourth hypotheses. Replicating prior  
8 work, binding values were positively correlated with stigmatization of victims, and individualizing  
9 values were positively correlated with sensitivity to victim suffering. The implicit causality object-  
10 bias for harm and force events was associated with explicitly less sensitivity to victim suffering,  
11 judgments of agents as less necessary and sufficient, greater explicit stigmatization of victims,  
12 and increased judgments that patients allowed, controlled and deserved harm and force events.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

MORAL VALUES AND CAUSAL ATTRIBUTION

**Table 5.** Spearman's rank-order correlations among moral values, judgments of victims as contaminated and injured, implicit causality object-bias for harm and filler events, and explicit causal judgments for agents and patients of harm/force and neutral filler events.

						IC object-bias		Explicit causal ratings			
		Binding	Individ.	Contam.	Injured	Harm	Filler	Harm: Agent	Harm: Patient	Fillers: Agent	Fillers: Patient
<b>Individ.</b>	Study 1	.047									
	Rep 1	.197									
	Rep 2	.139**									
<b>Contam.</b>	Study 1	.473***	-.101								
	Rep 1	.402***	-.086								
	Rep 2	.368***	.053								
<b>Injured</b>	Study 1	-.183***	.316***	-.325***							
	Rep 1	-.090	.284***	-.228***							
	Rep 2	.049	.160***	-.212***							
<b>IC Object-Bias</b>	<b>Harm</b>	Study 1	.311***	.000	.271***	-.097					
		Rep 1	.253***	-.063	.161	-.224***					
		Rep 2	.147***	-.080	.201***	-.161***					
	<b>Filler</b>	Study 1	.082	.001	-.024	.062	.206***				
		Rep 1	.097	.068	.008	-.007	.292***				
		Rep 2	.111	-.021	.088	-.036	.364***				
<b>Explicit causal ratings</b>	<b>Harm: Agent</b>	Study 1	-.227***	.164***	-.235***	.300***	-.296***	.079			
	<b>Harm: Patient</b>	Study 1	.271***	-.161	.312***	-.251***	.342***	-.042	-.515***		
	<b>Fillers: Agent</b>	Study 1	.021	.153	-.019	.107*	-.003	.126	.162	.027	
	<b>Fillers: Patient</b>	Study 1	.194***	.024	.166***	-.015	.140	.155	.020	.353***	.310***

Note. Study 1 (N = 459), Replication 1 (N = 249), Replication 2 (N = 788). Binding = Binding values. Individ. = Individualizing values. Contam. = Ratings of victims as contaminated. Injured = Ratings of victims as injured. Harm = IC object-bias: Implicit causality object-bias for events of harm and force. Fillers = IC object-bias: Implicit causality object-bias for neutral filler events. Harm: Agent = Causal contribution of agents for events of harm and force. Harm: Patient = Causal contribution of patients for events of harm and force. Fillers: Agent = Causal contribution of agents for neutral events. Fillers: Patient = Causal contribution of patients for neutral events. \*\*\* p <.001.

## MORAL VALUES AND CAUSAL ATTRIBUTION

**General Discussion**

*Moral Values and Causal Attribution.* Prior work found that moral values aimed at protecting groups and keeping people bound into tight-knit relationships – the “binding values” of loyalty, obedience to authority and purity (e.g., Graham et al., 2011) are associated with victim-blaming, stigmatizing judgments of victims as contaminated and tainted, and viewing victims as responsible (Niemi & Young, 2016). Because blame allocations were mediated by judgments of victim responsibility, it was speculated that binding values might involve a different understanding of the causal structure of events of harm and force. The results of the present research suggest that people who highly endorse binding values might interpret the causal structure of events slightly differently, as indicated by their responses to the implicit causality task, and by their explicit causal judgments. People who highly endorsed binding values are more likely to attribute causation to sentence objects over subjects in the implicit causality task across a range of harm and force events. In line with the interpretation that people high in binding values might interpret the causal structure of these events (e.g., “Bob coerced Amy”) differently, they rated sentence objects (patients) as more likely to have allowed, controlled, and deserved harm and force, and sentence subjects (agents) as less necessary and sufficient for harm and force.

Since immoral events in the case of binding values do not necessarily fit with the “moral dyad” in which an agent harmed a patient and the agent is causal whereas the affected patient is not (Gray, Young, Waytz, 2012), one possibility is that greater endorsement of binding values may lead people to spread causal responsibility across event participants – including surprising targets: victims (Niemi & Young, 2016). Another possibility that could be tested in future work is that people higher in binding values throw out the moral dyad altogether and abandon ascribing causation in the typical zero-sum, hydraulic manner (i.e., the more causal the victim is perceived, the less the perpetrator). Results from the measures of *explicit* causal judgments indicate however that people higher in binding values retain a zero-sum understanding of blame



## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 ascription. The more strongly people endorsed binding values, the more likely they were to  
4  
5 judge patients as having allowed, controlled and deserved harm, and the less likely they were to  
6  
7 judge agents of harm as necessary and sufficient. The more causal contribution people high in  
8  
9 binding values perceived in victims, the less they perceived in perpetrators.  
10

11  
12 By contrast, previous work found small associations between individualizing values and  
13  
14 increased judgments of perpetrators as causal contributors, and increased sensitivity to victims'  
15  
16 suffering (Niemi & Young, 2016). In the current research, similar associations were observed:  
17  
18 individualizing values correlated with explicit judgments of agents' contributions to harm and  
19  
20 increased sensitivity to victims' suffering (Table 5). Individualizing values were not associated  
21  
22 with selecting the subject (the harm-doer) in the implicit causality task. However, sensitivity to  
23  
24 victim suffering was related to selection of the harm-doer in the implicit causality task in two  
25  
26 studies (Table 5). This dissociation indicates an interesting area for further research:  
27  
28 individualizing values corresponded with sensitivity to suffering, yet only sensitivity to suffering  
29  
30 and not individualizing values predicted object-bias for harms in the IC.  
31

32  
33 It may be counterintuitive that binding values – *moral* values – motivate an  
34  
35 understanding of the causal structure of harm that places the causal source of the explanation  
36  
37 on the *affected* person. The link between binding values and object-bias may come about in at  
38  
39 least two ways. First, people high in binding values might be driven by relatively benign motives  
40  
41 – even though object-bias for harm correlated with rating victims as contaminated and patients  
42  
43 as more likely to have allowed, controlled, and deserved harmful events. We did not look at  
44  
45 participants' concern about recklessness, negligence, prudence or social harmony. It is possible  
46  
47 that people higher in binding values have increased concerns about victims (as imprudent  
48  
49 triggers of events that bring them harm) in keeping with their increased concern about group-  
50  
51 level order and other social moral concerns. Alternatively, the findings could represent a shift in  
52  
53 causal structure that comes along with moralized judgments that chastise a victim in order to  
54  
55 protect the status of the self and associates. This would be consistent with the purported  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 function of binding values to motivate behavior and attitudes that protect groups and the family  
4 unit before being concerned about the suffering of any one individual (e.g., Graham et al., 2011;  
5 Haidt, 2007), as well as prior findings of a positive association between binding values and  
6 status-seeking (Niemi & Young, 2013).  
7  
8  
9

10  
11 *Implicit Causality Task as a Social Psychology Tool.* Some researchers have proposed  
12 that implicit causality responses may differ systematically across different sorts of verbs  
13 because people draw on their experience with typical causes of those sorts of events (Bott &  
14 Solstad, 2014). However, evidence has been inconsistent, and the largest and most systematic  
15 investigations provide limited support for this claim (Ferstl et al., 2011; Pickering & Majid, 2007).  
16 Here, we find that individual variation in moral values is correlated with performance on the IC  
17 task specifically for highly morally relevant events. For these events, people higher in binding  
18 values were more likely to select the object as the cause, compared to people low in binding  
19 values. There was no consistent effect on morally neutral events. This raises a challenge for the  
20 IC literature, as current theories argue either that non-linguistic cognition is always relevant for  
21 IC or that it never is. Implicit verb causality is a nuanced phenomenon. Because our focus here  
22 is not theories of IC *per se*, we leave this as a challenge for future work.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

37 Future work will also need to examine individual differences in causal models based on  
38 people's experiences with these events, as well as whether people higher in binding values are  
39 more likely to select the object as causal in the implicit causality task because they have a  
40 broader temporal representation of harm and force events that presupposes a prior event in  
41 which the patient performed a bad action that made them deserving of "punishment" by the  
42 agent (Cf. Bott & Solstad, 2014). The finding that the event participants' gender mattered to  
43 implicit causality selections (i.e., people expect events involving women harming men to be  
44 better explained with reference to the object (men) compared to when men harm women)  
45 indicates that quick inferences about violence are jointly influenced by both stable moral values  
46 as well as gendered intuitions about deservingness for harm.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 The current research, therefore, reveals that the implicit verb causality task is  
4 psychologically informative based on which types of verbs are slotted into its simple “*Person A*  
5 *verbed Person B because ... he or she?*” format. Namely, people’s tendency to select the  
6 referent to the object in the IC task with “harm/force” verbs is reliably informative as to whether  
7 they will also consider victims of crimes “contaminated” or “injured”, and endorse “binding  
8 values”—moral values of loyalty, obedience to authority and preservation of purity—previously  
9 found to be associated with victim blaming (Niemi & Young, 2016). In contrast, we find that the  
10 implicit causality task applied to verbs with limited moral relevance, i.e., the neutral verbs, was  
11 not related to people’s moral values and social attitudes. In all likelihood, verb semantics drove  
12 participants’ responses in these cases, as prior work has demonstrated that linguistic structure  
13 explains implicit causality biases when participants’ individual differences aren’t taken into  
14 account (Hartshorne, 2013; Hartshorne & Snedeker, 2012). Thus, the IC task with the harm and  
15 force verbs was a useful social psychology tool and performance was distinct from performance  
16 in our explicit moral measures (i.e., moral values, ratings of victims), and from explicit causal  
17 judgments of agents as necessary and sufficient, and patients as having allowed, controlled or  
18 deserved the events.

19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37 *The Affordances of Multiple Measures of Causality.* Regarding the explicit measures of  
38 causality, the current work is the first to link implicit causality selections with people’s judgments  
39 of agents’ necessity and sufficiency, and patients’ allowing, controlling and desert of events.  
40 Other researchers have examined how verbs’ implicit causal biases vary with other kinds of  
41 causally-relevant information about people (e.g., covariation (Brown & Fish, 1983; Rudolph,  
42 2008)). Our aim goes beyond linking implicit causality behavior and explicit causal judgments:  
43 we examine the relationship between these measures and both individual differences in stable  
44 moral values and moral judgments of situations. Nevertheless, we chose to assess agents’  
45 necessity and sufficiency because these are typically considered conditions relevant to being  
46 the cause. To make sure that the way in which we ask these questions is not confounded with  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 explicit moral judgments, we asked “Would [Amy] have been [coerced] by someone else?” to  
4 assess necessity, and “Would [Bob] [coerce] someone else?” to assess sufficiency. One might  
5 argue that in order to answer these questions, participants still have to ask themselves whether  
6 Amy is a pushover, or whether Bob is manipulative. We could also have asked more directly:  
7 “Who caused the rape? Amy or Bob?” It’s likely that such direct questioning would alert  
8 participants to social desirability concerns or trigger reactive affect about victim-blaming.

9  
10  
11  
12  
13  
14  
15  
16 Measuring participants’ judgments of agents’ and patients’ explicit causal contributions  
17 more covertly with multiple questions not only helped circumvent social desirability concerns. It  
18 also enabled participants to assign their judgments more freely. Instead of using a bipolar scale  
19 such as “agent-caused vs. patient-caused”, which would require participants to treat causation  
20 in a zero-sum manner across the dyad, our items measuring agents’ necessity and sufficiency  
21 and patients’ capacity to control, allow, and deserve the event let us determine whether  
22 participants treated agents’ and patients’ contributions in a zero-sum manner even when these  
23 were measured unconstrained. We found that participants do indeed treat agents’ and patients’  
24 explicit contributions as though they are hydraulically related (i.e., when agents are rated more  
25 causal, patients are rated less causal). In addition, their explicit responses correlated with  
26 implicit causality responses, which are bipolar in nature. Finally, the use of multiple explicit  
27 causality items with scaled response options revealed that higher endorsement of binding  
28 values is not just associated with broad over-attribution of causal responsibility to agents *and*  
29 patients of harm – binding values are associated specifically with over-attribution to patients.

30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
Inquiring about participants’ explicit moral judgments, implicit causal selections, and  
explicit moral values and judgments also allowed us to observe whether and how these  
variables interrelated. Most notably, because of the potential consequences for harmed people,  
binding values of loyalty, obedience to authority, and preservation of purity were related to  
stigmatization of victims (replicating previous findings, Niemi & Young, 2016), increased explicit  
causal judgments of patients and reduced causal judgment of agents, and implicit causality

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 object-bias for harm and force. No relationship was observed between binding values and  
4 sensitivity to victims' suffering (ratings of victims as injured), whereas implicit causality object-  
5 bias for harm and force *did* correlate with reduced sensitivity for victim suffering in two studies.  
6  
7 This finding suggests two potential sources driving object-bias for harm in the implicit causality  
8 task: (1) callousness, and (2) binding values.  
9  
10  
11  
12

13  
14 This research demonstrates the advantages of measures that tap multiples levels of  
15 awareness and, in particular, the advantages of the implicit causality task as a measure of  
16 people's intuitions about causation in the case of harm and force. Since the task is repeated  
17 over several trials, the experimenter can embed numerous foils, including positive and neutral  
18 events. As the response options are limited to just "he" or "she", people are likely to  
19 underestimate the extent to which any individual choice may be informative. Yet, in fact, we find  
20 that stripping a range of events involving harm and force down to their most minimal possible  
21 descriptions (e.g., "Bob coerced Amy because") and determining the likelihood that participants  
22 select the object as referent results in an informative measure about morality. Most reliably, it  
23 informs about their tendencies toward victim stigmatization and their moral commitments: their  
24 valuation of loyalty, obedience to authority, and concern about preservation of purity.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

37 These latter social-moral attitudes are attitudes that those in military, legal, and clinical  
38 settings, who lead, litigate, and care for harmed people might prefer to guard or conceal. Thus,  
39 there is viable research utility for the implicit causality task, for example, in testing its use as a  
40 covert measure of attitudes toward stigmatized populations, e.g., rape victims, minorities in  
41 various settings. More broadly, in many organizational settings, it is important that attributions of  
42 causation can be measured covertly in the service of understanding how people think about  
43 blame and responsibility, and this task shows promise for that purpose.  
44  
45  
46  
47  
48  
49  
50

51 **Conclusion.** In sum, prior work showed that "binding values" of loyalty, obedience to  
52 authority and purity are reliably linked with stigmatization and blame of victims whereas  
53 "individualizing values", concerned with broadly preventing harm and unfairness, are linked with  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 sensitivity to victim suffering and attributing responsibility to perpetrators. The current work  
4 suggests a mechanism for these associations in judgments of agents and patients as causal  
5 contributors. Moral values are directly tied to judgments of agents as necessary and sufficient  
6 contributors, and patients as having allowed, controlled and deserved the outcome. Moreover,  
7 this work indicates that moral values predict a shift in people's expectations about who caused a  
8 range of events of harm and force: the person who did it, or the person who had it done to them.  
9 Taken together, the results indicate that people with different moral values might differ in their  
10 initial assessments of the causal source of harmful events, which may, in turn, relate to their  
11 explicit attitudes about stigmatization and blame.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 **Author Contributions.** All authors contributed to the study design, data interpretation,  
25 manuscript preparation and revisions, and approved the final version of the manuscript for  
26 submission. Testing, data collection and analysis were performed by BLINDED and BLINDED.  
27  
28  
29  
30  
31  
32

33 **Competing Interests.** The authors declare no competing interests.  
34  
35  
36

37 **Acknowledgments.** We gratefully acknowledge the helpful comments and suggestions of Fiery  
38 Cushman, Felipe De Brigard, John Doris, Joshua Greene, Thomas Icard, Joshua Knobe, Jorie  
39 Koster-Hale, Laurie Paul, Jonathan Phillips, Steven Pinker, Steve Sloman, Jesse Snedeker,  
40 CUNY 2016 attendees, the Morality Lab at BC, the Events in Language Reading Groups at  
41 Harvard and MIT, and the Moral Psychology Research Group, and funding by a Sage Young  
42 Scholar Award.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MORAL VALUES AND CAUSAL ATTRIBUTION

## References

- Alicke, M. D. (2000). Culpable control and the psychology of blame. *Psychological Bulletin*, *126*, 556-574.
- Alicke, M. D., Mandel, D. R., Hilton, D. J., Gerstenberg, T., & Lagnado, D. A. (2015). Causal conceptions in social explanation and moral evaluation: A historical tour. *Perspectives in Psychological Science*, *10*, 790-812.
- Arnold, J. E. (2015). Women and men have different discourse biases for pronoun interpretation. *Discourse Processes*, *52*, 77-110.
- Bates, D., Maechler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effect models using lme4. *Journal of Statistical Software*, *67*, 1-48.
- Bott, O., & Solstad, T. (2014). From verbs to discourse: A novel account of implicit causality. In B. Hemforth, B. Mertins & C. Fabricius-Hansen (Eds.), *Psycholinguistic approaches to meaning and understanding across languages*. (pp. 213-251). Cham, Switzerland: Springer.
- Brown, R. & Fish, D. (1983). The psychological causality implicit in language. *Cognition*, *14*, 237-273.
- Cushman, F. (2008). Crime and punishment: Distinguishing the roles of causal and intentional analyses in moral judgment. *Cognition*, *108*, 353-380.
- Ferstl, E. C., Garnham, A., & Manouilidou, C. (2011). Implicit causality bias in English: A corpus of 300 verbs. *Behavior Research Methods*, *43*, 124-135.
- Garvey, C., & Caramazza, A. (1974). Implicit causality in verbs. *Linguistic Inquiry*, *5*, 459-464.
- Glick, P., & Fiske, S. (1996). The ambivalent sexism inventory: Differentiating hostile and benevolent sexism. *Journal of Personality and Social Psychology*, *70*, 491-512.
- Graham, J., Nosek, B. A., Haidt, J., Iyer, R., Koleva, S., & Ditto, P. H. (2011). Mapping the moral domain. *Journal of Personality and Social Psychology*, *101*, 366-385.

## MORAL VALUES AND CAUSAL ATTRIBUTION

- 1  
2  
3 Gray, K., Young, L., & Waytz, A. (2012). Mind perception is the essence of morality.  
4  
5 *Psychological Inquiry*, 23, 101-124.  
6  
7 Haidt, J. (2007). The new synthesis in moral psychology. *Science*, 316, 998-1002.  
8  
9 Hartshorne, J. K. (2013). What is implicit causality? *Language, Cognition and Neuroscience*, 29,  
10  
11 804-824.  
12  
13 Hartshorne, J. K., & Snedeker, J. (2012). Verb argument structure predicts implicit causality:  
14  
15 The advantages of finer-grained semantics. *Language and Cognitive Processes*, 28,  
16  
17 1474-1508.  
18  
19 Hartshorne, J. K., Sudo, Y., Uruwasi, M. (2013). Are implicit causality pronoun resolution biases  
20  
21 consistent across languages and cultures? *Experimental Psychology*, 60, 179-196.  
22  
23 Heider, F. (1958). *The psychology of interpersonal relationships*. New York: John Wiley & Sons,  
24  
25 Inc.  
26  
27 Hesslow, G. (1988). The problem of causal selection. In D. J. Hilton (Ed.), *Contemporary*  
28  
29 *science and natural explanation: Commonsense conceptions of causality* (pp. 11-32).  
30  
31 Brighton, England: Harvester Press.  
32  
33 Hilton, D. J. (1990). Conversational processes and causal explanation. *Psychological Bulletin*,  
34  
35 107, 65-81.  
36  
37 Holtgraves, T. M. (2002). *Language as social action: Social psychology and language use*. New  
38  
39 Jersey: Lawrence Erlbaum Associates.  
40  
41 Kipper-Schuler, K. (2006). *VerbNet: A broad-coverage, comprehensive verb lexicon*. Ph.D.  
42  
43 thesis, University of Pennsylvania.  
44  
45 Levin, B. (1993). *English verb classes and alternations*. Chicago: Chicago University Press.  
46  
47 Lombrozo, T. (2006). The structure and function of explanations. *Trends in Cognitive Sciences*,  
48  
49 10, 464-470.  
50  
51 Lombrozo, T. (2009). The role of moral commitments in moral judgment. *Cognitive Science*, 33,  
52  
53 273-286.  
54  
55  
56  
57  
58  
59  
60



## MORAL VALUES AND CAUSAL ATTRIBUTION

1  
2  
3 Malle, B. F., Guglielmo, S., & Monroe, A. E. (2014). A theory of blame. *Psychological Inquiry*,  
4  
5 25, 147-186.

6  
7 Nappa, R., & Arnold, J. E. (2014). The road to understanding is paved with the speaker's  
8  
9 intentions: Cues to the speaker's attention and intentions affect pronoun  
10  
11 comprehension. *Cognitive Psychology*, 70, 58-81.

12  
13 Niemi, L. & Young, L. (2016). When and why we see victims as responsible: The impact of  
14  
15 ideology on attitudes toward victims. *Personality and Social Psychology Bulletin*, 42,  
16  
17 1227-1242.

18  
19 Niemi, L & Young, L. (2013). Caring across boundaries versus keeping boundaries intact: Links  
20  
21 between moral values and interpersonal orientations. *PLoS One*, 8(12), e81605.

22  
23 Pickering, M. J., & Majid, A. (2007). What are implicit causality and consequentiality? *Language*  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
*and Cognitive Processes*, 22, 780-788.

Rudolph, U. (2008). Covariation, causality, and language: Developing a causal structure of the  
social world. *Social Psychology*, 39, 174-181.

Rudolph, U., & Forsterling, F. (1997). The psychological causality implicit in verbs: A review.  
*Psychological Bulletin*, 121, 192-218.

Shaver, K. G. (1985). *The attribution of blame: Causality, responsibility, and*  
*blameworthiness*. New York: Springer-Verlag.

## MORAL VALUES AND CAUSAL ATTRIBUTION

## Footnotes

<sup>1</sup> Attention checks failure involved responding with 1 or 2 on a Likert-scale of agreement with the item "It is better to do good than bad," or 5 or 6 on the scale measuring how relevant it was to their criteria of right or wrong: "Whether or not someone was good at math" from the Moral Foundations Questionnaire (the standard "attention check items" from the MFQ), and completion of any of four blocks of MFQ questions in under 10 s. Original sample = 648, any one error was sufficient for exclusion. We excluded 189 additional individuals who failed attention checks.

<sup>2</sup> Attention checks failure criteria for Replication Datasets 1 and 2 were identical to Study 1. Replication Dataset 1 attention check failures = 135; exclusions for failing to complete the study = 315, likely because following the IC portions, this study involved a lengthy pilot. Replication Dataset 2 attention check failures = 284.

<sup>3</sup> Two additional verbs, "confused" and "punished," were omitted from analyses over concern about the neutrality of *punished*, and to balance the representation of typically subject- and object-biased verbs with its removal.

<sup>4</sup> Similar statistical modeling approaches have been implemented in psycholinguistics research (e.g., Nappa & Arnold, 2014). We did not include random slopes because many models failed to converge when random slopes were included.

<sup>5</sup> Harms used in prior work conducted by Niemi and Young (2016).

For Peer Review

1  
2  
3 **Supplementary Materials for Moral values in causal attribution:**  
4 **Evidence from the implicit verb causality task and explicit judgments**  
5  
6  
7  
8

9 **A: Notes on Methodological Differences between Study 1 and the Replication Datasets**  
10

11 **B: Gender Condition and Implicit Causality Object-Bias from Study 1 and the Replication**  
12 **Datasets**  
13  
14

15 **C: Demographic Controls and Implicit Causality Object-Bias**  
16  
17

18 **D: Individualizing values and Implicit Causality Object-Bias**  
19  
20  
21  
22  
23  
24

25 **A: Notes on Methodological Differences between Study 1 and the Replication Datasets**  
26

27  
28 *Replication Dataset 1.* As in Study 1, participants completed a block of the implicit causality task  
29 and then entered a separate block; this time the events were re-presented with the pronoun  
30 they had selected in the first block (e.g., “Max verbed Jess because he...”) and an empty text  
31 box appeared after the pronoun with the prompt: “please finish the sentence.” Participants typed  
32 a completion to each sentence. They also filled out measures of demographics, victim  
33 stigmatization and sensitivity, and moral values as in the previous study (data and materials  
34 available at [https://github.com/ BLINDED FORREVIEW](https://github.com/BLINDED FORREVIEW)).  
35  
36  
37  
38  
39  
40  
41  
42

43 *Replication Dataset 2.* As in the previous studies, participants completed the implicit causality  
44 task, but the task here involved a total of 48 verbs. See Table 1 in the main text for the complete  
45 list of verbs. Participants also filled out measures of demographics, victim stigmatization and  
46 sensitivity, and moral values as in the previous studies. They completed the Ambivalent Sexism  
47 Inventory (Glick & Fiske, 1996); the present analyses did not involve the ASI. Data and  
48 materials are available at [https://github.com/ BLINDED FORREVIEW](https://github.com/BLINDED FORREVIEW).  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## B: Gender Condition and Implicit Causality Object-Bias from Study 1 and the Replication

### Datasets

To investigate whether the gender condition (male-verbed-female *versus* female-verbed-male) predicted the implicit causality object-bias for harm/force verbs, we first computed a generalized linear mixed-effects regression model in which verb type (harm/force (coded as 0) *versus* neutral filler (coded as 1)) and gender condition (male-verbed-female (coded as 0) *versus* female-verbed-male (coded as 1)) were included as fixed predictors of the propensity to select the object (coded as 1) *relative to* the subject (coded as 0) as the referent (participant and verb were both included as random effects with random intercepts only). There was a significant interaction between verb type and gender condition in Study 1 and in the Replication Datasets (see **Supplementary Table 1**). To further interrogate these significant interaction effects, generalized linear mixed-effects models were computed for harm/force verbs and neutral filler verbs, taken separately.

**Supplementary Table 1.** *The results of two generalized linear mixed-effects regression models—each with verb type and gender condition as predictors of selecting the object versus the subject as the referent.*

	<i>b</i>	<i>SE</i>	<i>Z</i>	<i>p</i>	95% CI
<b>Study 1</b>					
Verb Type	1.67	.40	4.19	< .0001	[.89, 2.45]
Gender Condition	.53	.11	4.89	< .0001	[.32, .75]
Verb Type x Gender Condition	-.99	.10	-10.29	< .0001	[-1.18, -.80]
<b>Replication Dataset 1</b>					
Verb Type	1.52	.29	3.95	< .0001	[.77, 2.28]
Gender Condition	.71	.15	4.77	< .0001	[.42, 1.00]
Verb Type x Gender Condition	-.78	.13	-5.99	< .0001	[-1.03, -.52]
<b>Replication Dataset 2</b>					
Verb Type	1.45	.30	4.81	< .0001	[.86, 2.04]
Gender Condition	.62	.08	8.28	< .0001	[.48, .77]
Verb Type x Gender Condition	-.87	.05	-17.74	< .0001	[-.97, -.78]

*Note.* Main Experiment ( $N = 459$ ), Replication Dataset ( $N = 249$ ), Replication Dataset 2 ( $N = 788$ ). All 95% CIs are for the beta-estimates.

For harm/force verbs, a generalized linear mixed-effects regression model was computed for which gender condition was included as the fixed predictor of the propensity to select the

1  
2  
3 object (coded as 1) *relative to* the subject (coded as 0) as the referent (participant and verb were  
4 both included as random effects with random intercepts only). This analysis yielded a significant  
5 effect of gender condition on the likelihood of selecting the object versus the subject as the  
6 referent in Study 1 ( $b = .63$ ,  $SE = .13$ ,  $Z = 4.98$ ,  $p < .0001$ , 95% CI = [.38, .87]). We obtained the  
7 same pattern of results in Replication Dataset 1 ( $b = .80$ ,  $SE = .18$ ,  $Z = 4.35$ ,  $p < .0001$ , 95% CI =  
8 [.44, 1.16]) and Replication Dataset 2 ( $b = .69$ ,  $SE = .08$ ,  $Z = 8.31$ ,  $p = .0001$ , 95% CI = [.53, .85]).  
9  
10 In all datasets, when women harmed/forced men, participants were more likely to select the object  
11 than the subject as the referent (i.e., there was a more pronounced object-bias).  
12  
13  
14  
15  
16  
17  
18  
19  
20

21 For neutral filler verbs, there was a significant effect of gender condition on the  
22 propensity for selecting the object over the subject as the referent in the main experiment ( $b = -$   
23  $.43$ ,  $SE = .11$ ,  $Z = -3.80$ ,  $p = .0001$ , 95% CI = [-.66, -.21]); there was no effect in Replication  
24 Dataset 1 ( $b = -.04$ ,  $SE = .14$ ,  $Z = -.28$ ,  $p = .78$ , 95% CI = [-.32, .24]); the effect returned in  
25 Replication Dataset 2 ( $b = -.24$ ,  $SE = .07$ ,  $Z = -3.47$ ,  $p = .0005$ , 95% CI = [-.38, -.11]). The  
26 significant effects for neutral filler verbs were in the opposite direction of the effects for  
27 harm/force verbs. To sum up, moral values aside, participants generally were more likely to  
28 select men for harm/force events in the implicit causality task regardless of whether they were  
29 the subject (the “perpetrator” of harm/force) or the object (the “victim” of harm/force).  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 However, binding values withstood this gender effect. Importantly, when gender  
41 condition was added to the generalized linear mixed-effects models that already included  
42 binding values as a predictor of the propensity to select the object relative to the subject as the  
43 referent, participants higher in binding values remained significantly more likely to select the  
44 object over the subject as referent for harm/force verbs in Study 1 ( $b = .41$ ,  $SE = .06$ ,  $Z = 6.76$ ,  $p$   
45  $< .0001$ , 95% CI = [.29, .52]), Replication Study 1 ( $b = .50$ ,  $SE = .10$ ,  $Z = 4.95$ ,  $p < .0001$ , 95%  
46 CI = [.30, .69]), and Replication Study 2 ( $b = .21$ ,  $SE = .05$ ,  $Z = 4.65$ ,  $p < .0001$ , 95% CI = [.12,  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60 .30]).

### C: Demographic Controls and Implicit Causality Object-Bias

Having found that for harm/force verbs, binding values significantly predict the likelihood of selecting the object over the subject as referent (“object-bias”), we next expanded the generalized linear mixed-effects regression models to include political orientation, gender, and religiosity as additional fixed predictors along with binding values. Given that prior work has identified relationships between binding values and political orientation, gender, and religiosity (Graham et al., 2011), we wanted to ensure that binding values predicted the implicit causality object-bias above and beyond these other variables. More specifically, binding values, political orientation, gender (0 = male, 1 = female), and religiosity were included as fixed predictors of the propensity to select the object (coded as 1) *relative to* the subject (coded as 0) as the referent for the harm/force verbs. Full results for these models for Study 1 and the replication datasets are depicted in **Supplementary Table 2**. Binding values remained a consistent, significant predictor of the object-bias in all three datasets after statistically controlling for political orientation, gender, and religiosity. An inconsistent effect of gender appeared in two cases in Study 1 and Replication Dataset 2 such that male participants were more likely to exhibit object-bias than female participants.

**Supplementary Table 2.** *The results of two generalized linear mixed-effects regression models—each with binding values, political orientation, gender, and religiosity as predictors of the propensity to select the object as the referent for harm and force events.*

	<i>b</i>	<i>SE</i>	<i>Z</i>	<i>p</i>	95% CI
<b>Study 1</b>					
Binding Values	.21	.08	2.73	.006	[.06, .36]
Political Orientation	-.07	.04	-1.59	.111	[-.15, .01]
Gender	-.58	.12	-4.92	< .0001	[-.82, -.35]
Religiosity	.09	.03	2.69	.007	[.02, .15]

**Replication****Dataset 1**

Binding Values	.53	.13	4.08	< .0001	[.28, .79]
Political Orientation	.08	.06	1.31	.190	[-.04, .21]
Gender	-.35	.18	-1.92	.055	[-.71, .01]
Religiosity	.00	.05	-.08	.938	[-.11, .10]

**Replication****Dataset 2**

Binding Values	.24	.06	4.16	< .0001	[.13, .35]
Political Orientation	-.02	.03	-.60	.55	[-.08, .04]
Gender	-.19	.09	-2.09	.037	[-.36, -.01]
Religiosity	-.04	.02	-1.54	.12	[-.08, .01]

Note. Study 1 ( $N = 459$ ), Replication Dataset 1 ( $N = 249$ ), Replication Dataset 2 ( $N = 788$ ). All 95% CIs are for the beta-estimates.

**D: Individualizing values and Implicit Causality Object-Bias**

We investigated whether individualizing values also predict a subject- or object-bias in the implicit causality task. Previously, individualizing values were found to be positively associated with perpetrator blame (Niemi & Young, 2016). This association was notably weaker than the associations between binding values and judgments of victims as blameworthy and responsible. Therefore, we did not have strong expectations regarding the implicit causality behavior of participants high in individualizing values. Nevertheless, increased selection of the subject for harm/force verbs would be consistent with the prior findings of increased perpetrator blame. To investigate this, a generalized linear mixed-effects regression model was computed in which verb type (harm/force (coded as 0) *versus* neutral filler (coded as 1)) and individualizing values were included as fixed predictors of the propensity to select the object (coded as 1) *relative to* the subject (coded as 0) as the referent.

In Study 1, there was no significant main effect of individualizing values ( $b = -.02$ ,  $SE = .08$ ,  $Z = -.27$ ,  $p > .05$ , 95% CI = [-.18, .13]), no significant main effect of verb type ( $b = .88$ ,  $SE =$

1  
2  
3 .52,  $Z = 1.72$ ,  $p > .05$ , 95% CI = [-.13, 1.90]), and no significant interaction ( $b = .06$ ,  $SE = .07$ ,  $Z$   
4 = .85,  $p > .05$ , 95% CI = [-.08, .20]). There was a small but significant interaction in Replication  
5 Dataset 1 ( $b = .23$ ,  $SE = .09$ ,  $Z = 2.44$ ,  $p = .01$ , 95% CI = [.04, .41]), but there were no  
6 significant main effects ( $ps > .05$ ). Despite this significant interaction effect, there was still no  
7 effect of individualizing values on the propensity to select the object *relative to* the subject as the  
8 referent for the subset of harm/force verbs ( $p > .05$ ) or the subset of neutral filler verbs ( $p > .05$ ),  
9 when these verb types were modeled separately. In Replication Dataset 2, there was no  
10 significant main effect of individualizing values ( $b = -.08$ ,  $SE = .06$ ,  $Z = -1.49$ ,  $p > .05$ , 95% CI =  
11 [-.19, .03]) and no significant interaction between individualizing values and verb type ( $b = .07$ ,  
12  $SE = .04$ ,  $Z = 1.81$ ,  $p > .05$ , 95% CI = [-.01, .14]). There was, however, a small but significant  
13 main effect of verb type ( $b = .72$ ,  $SE = .35$ ,  $Z = 2.06$ ,  $p = .04$ , 95% CI = [.04, 1.40]). Therefore,  
14 across the three datasets, it is reasonable to conclude that binding values, not individualizing  
15 values, are associated with the object-bias (for harm/force verbs and not neutral filler verbs).  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60