

Social Psychology

Asymmetries in Attributions of Blame and Praise, Intent, and Causality: Free Will, Responsibility, and the Side-effect Effect

Adrien Fillon¹, Subramanya P. Chandrashekar², Gilad Feldman³^a

¹ Science and Innovation Policy & Studies, University of Cyprus, Nicosia, Cyprus, ² Norwegian University of Science and Technology, Trondheim, Norway, ³ Department of Psychology, University of Hong Kong, Hong Kong SAR

Keywords: free will, experimental philosophy, attributions, side-effect effect, blame, praise

<https://doi.org/10.1525/collabra.128423>

Collabra: Psychology

Vol. 11, Issue 1, 2025

The Side-Effect Effect (SEE) is the phenomenon that negative side-effects elicit stronger attributions of intent and blame than intent and praise for positive side-effects. There are similar documented asymmetries showing stronger free will attributions to negative than to positive, and stronger associations between free will attributions and blame for negative outcomes than associations between free will attributions and praise for positive outcomes. Together, these are two well-known paradigms in experimental philosophy that have thus far mostly been studied separately. Given that they both examine similar domains regarding agency, intent, and responsibility, we aimed to integrate the two paradigms to examine possible joint effects and interactions. We used the classic SEE scenario with within and between designs, manipulated free will by contrasting deterministic versus indeterministic universes, and measured free will attributions. In two experiments (overall $N = 1520$), we found support for side-effect effects regarding attributions of intentionality and knowledge (Study 1: $d = 0.58-1.77$; Study 2: $d = 0.61-1.75$). We found a strong association between blame/praise and free will attributions, even when controlling for intent and knowledge. Finally, we found that when participants were asked to imagine a counterfactual and report praise or blame based on the experimental condition, blame was more strongly attributed to hypothetical harmful outcomes than praise to helpful outcomes. We found no consistent support for an interaction between the two paradigms, suggesting that they uniquely affect attributions. All materials, data, and code are available on: <https://osf.io/z3g6d/>

In the last decade, there has been increasing interest in moral social cognition, examining how people perceive, interpret, and understand moral behavior. Experimental philosophy has brought philosophy into the lab, testing lay beliefs and folk psychology of abstract philosophical questions. This work has led to interesting observations revealing cognitive processes regarding the way that people think regarding philosophical domains such as intent, morality, and free will. In the present investigation, we set out to combine two of the most well-known paradigms in experimental philosophy – the classic side-effect effect (SEE) impacting attributions of intent, and the classic thought experiments regarding an (in)deterministic universe impacting attributions of free will and moral responsibility. Our goal was to investigate the interplay between the SEE and free will attribution paradigms.

Side-Effect Effect

SEE is the phenomenon that harmful outcomes of an action are perceived as more intentional than helpful outcomes, even when the agent had no particular desire to bring about these outcomes (Knobe, 2003). Studies of the phenomenon typically introduce participants to the following vignette (brackets describe the manipulation):

The vice-president of a company went to the chairman of the board and said, ‘We are thinking of starting a new program. It will help us increase profits, but it will also [positive condition - **help**; negative condition - **harm**] the environment.’

The chairman of the board answered, ‘I don’t care at all about [positive condition - **helping**; negative condition - **harming**] the environment. I just want to make as much profit as I can. Let’s start the new program.’

^a Corresponding author:

Gilad Feldman, Department of Psychology, University of Hong Kong, Hong Kong SAR; gfeldman@hku.hk

They started the new program. Sure enough, the environment was [positive condition - **helped**; negative condition - **harmed**].

Now consider a seemingly straightforward question: Did the chairman of the board [positive condition - **help**; negative condition - **harm**] the environment intentionally?

In the original experiment, 82% of participants in the negative outcome condition reported that the agent intentionally brought about the side-effect. In contrast, only 33% of participants in the positive outcome condition attributed intentionality to the agent described in the scenario.

People attributed more intentionality and blame to negative side-effects than intentionality and praise to positive side-effects, although the descriptions were identical, aside from the outcomes. Thus, SEE exemplifies a blame-praise judgments asymmetry (Hindriks, 2008) and its links to folk perceptions of intentionality (Chandrashekar et al., 2022; Malle & Knobe, 1997). The extant research has also reported SEE regarding attributions of knowledge (Beebe & Buckwalter, 2010; Beebe & Jensen, 2012). The SEE was proposed as an alternative account to the earlier view that the motivations and intent of the agent are the basis of the intentionality attributed to a behavior (Heider, 1958; Ohtsubo, 2007; Shultz & Wells, 1985). The SEE account proposed that outcomes influence the perceiver's reasoning about the intentionality of the described behavior.

Judgments of blame (vs. praise) are affected by multiple sources of information related to the outcome, including the agent's foreseeability of the outcome, the intent of the agent causing the harm, and counterfactuals about the agent's action (Cushman et al., 2008; Laurent et al., 2019). Much of the impetus in explaining the SEE has been focused on the intuitions of intentionality.

Since the first demonstration (Knobe, 2003), SEE of intentionality has been considered a fairly robust effect (Beebe & Buckwalter, 2010; Cova & Naar, 2012; Feltz, 2007; Guglielmo & Malle, 2010; Klein et al., 2018; Laurent et al., 2019), and subsequently has been documented in other aspects such as causality (Tannenbaum et al., 2007), desire (Pettit & Knobe, 2009), and action versus in-action (Cushman et al., 2008). In addition, further work explaining the underlying cognitive processes that bring about asymmetry in the intuitions of ordinary subjects related to the SEE notes the role of emotions (Zucchelli et al., 2019) and individuals' personality differences (Cokely & Feltz, 2009).

Judgments of moral responsibility take into account several different aspects such as causality, intent, and counterfactuals about what could have been different (Malle, 2021). Blame serves the function of regulating behaviors of individuals in a society that promotes adherence to a set of moral standards (Monroe & Malle, 2019; Tetlock et al., 2010). Moreover, blame as an aspect of regulation extends to unintentional outcomes. For example, Monroe and Malle (2019) found that blame is constrained by the evidence that one's moral judgment is justified. Komatsu et al. (2021) found that robots are blamed more for inaction than humans when they fail to save lives because people think ro-

bots can prevent death better than humans. In other words, blame judgments also take into account the preventability, and the possibility that an agent could have taken steps to prevent an adverse outcome modulates the assessment of blame (Martin et al., 2019; Weiner, 1995).

On the other side, praise has often been overlooked. Indeed, while both are judgments regarding intentionality and causality, praise appears less sensitive to these features, and more in line with general features about an individual's stable, underlying character traits (Anderson et al., 2020). Blame seems to be more about the action, whereas praise seems to be more about the person who performed the action. This may explain why studies have documented blame and praise asymmetries in that they elicit very different attributions.

Linking Free Will and SEE: Free Will and Intent Attributions to Side-effects

Free will is often understood as a necessary condition for moral responsibility, because people perceive accountability as dependent on the person's capacity to have chosen to do otherwise (Monroe et al., 2014). In other words, that the agent has chosen their behavior freely, which may suggest stronger responsibility for his/her own actions. Empirical studies found support for the view that negative actions and outcomes were attributed stronger free will than positive ones, even for non-moral scenarios (Feldman et al., 2016; Fillon et al., 2022; Genschow & Vehlow, 2021). We therefore speculated that an agent in a situation involving a harmful outcome scenario is attributed more free will than an agent with a beneficial outcome, even when the outcome was a side-effect.

The side-effect effect paradigm has been used to demonstrate asymmetries in the attribution of intent and blame/praise to seemingly unintended side-effects. The attribution of intent is therefore not only one of the factors associated with the attribution of blame and praise, but, looking at the reverse causal chain, intent is affected by the attribution of blame, so that the need for blame and holding someone accountable leads to stronger attributions of intent (Malle et al., 2022; Monroe & Malle, 2019). Therefore, if even unintended harmful side-effects elicit higher intent, then it is possible that the "bad is freer than good" paradigm identified for free will attributions (Feldman et al., 2016) also extends to lesser chosen or "free" side-effects. That is, if the intentionality side-effect effect extends to free will attributions, then even if the protagonist (e.g., the chairman) only chooses to do something because of focusing on a different unrelated reason (e.g., to increase profits) and that this choice is driven by external pressures (e.g., the board and the shareholders), then with harmful outcomes (e.g., environment is harmed) the protagonist is still attributed as having more free will and the capacity for choice to do otherwise.

Further, examining intent and free will attributions together also helps make clearer the differences between them and their possible links in theories of blame and blame models. For example, in Malle et al. (2014)'s Theory of Blame they provided a "Path Model of Blame" with many

different factors, including “intentionality” (“whether the agent brought about the event intentionally”) and “capacity” (“whether the agent could have prevented the event”), yet missing the component of “free will” or “choice” (whether the agent could have chosen whether to prevent the event or not). Choice is loosely related to some of the other factors in the path model, such as “obligation” which serves as an external pressure limiting choice, yet goes far beyond that in capturing internal and external factors that may have restricted choice (Feldman, 2017). Studying intentionality and free will together by first using two experimental philosophy paradigms that focus on free will and intentionality, and then measuring both free will, intentionality, and blame attributions, can help shed light on 1) the associations between the three and 2) how each factor is affected by manipulations that impact free will and intentionality.

Linking Free Will and SEE: Manipulation of Both Agency and Outcome Valence

Experimental philosophy used thought experiments to provide additional insights regarding causality, determinism, and compatibilism (Feltz & Cova, 2014; Nahmias et al., 2007; Nichols & Knobe, 2007). We adapted this methodology to combine the thought experiments used to study SEE and free will. As the background context for the classic experimental philosophy SEE chairman scenario, we added the classic experimental philosophy manipulation of free will by describing the universe as either being deterministic or indeterministic.

The side-effect effect was initially demonstrated about intentionality: that when negative outcomes occur, to hold people accountable people judge other’s behavior linked to that outcome to be an intended consequence of their action. Intentionality attributions are evaluations of whether an outcome of an associated action was planned. For example, when the protagonist has no foreknowledge of the negative outcome, then intentionality attributions are weaker (Laurent et al., 2019).

Free will attributions are focused on agency and choice: whether people are perceived to have had the choice to do otherwise, without internal or external constraints (Feldman, 2017).

While both attributions are associated with blame (Malle et al., 2014), they are conceptually and empirically different (Feldman et al., 2016; Phillips & Knobe, 2009). One important difference is that free will is (mostly) the capacity for action regardless of constraints, both internal and external, and regardless of the outcome, whereas intentionality is a purely internal process, and focused on an association with an outcome. However, these nuanced differences in peoples’ understanding of free will and intention have so far not been comprehensively examined in the literature (Feldman, 2017).

In the present investigation, we measured intention and free will attributions, and we manipulated SEE and free will environment to assess how intentionality and free will attributions covary, and how they might be differentiated when judging unintentional harm and help outcomes. We

considered the control condition for the manipulation of free will universe as a replication of the side-effect effect.

Extension: Attributions of Regret and Moral Responsibility

We also added an extension aiming to examine attributions of regret in the context of free will and SEE. Fillon et al. (2022) examined the relationship between agency and regret, and reported stronger regret attributed to exceptionality compared to regret, and with stronger regret attributed in an in-deterministic universe compared to deterministic universe, with no support for an interaction. Additionally, they reported that regret attributions were positively associated with free will and moral responsibility attributions ($r = 0.20 - 0.42$). Their findings overall suggested that when things go badly, a stronger sense of agency is related to feeling more responsible for the negative outcome and feeling stronger regret for it. Given that in the SEE there is a manipulation of outcome valence, we were interested whether: 1) we would be able to replicate the pattern of results for regret, and 2) whether this pattern of stronger responsibility would also translate to taking credit for positive outcomes, and whether that would be impacted by manipulation of agency.

The Present Investigation

We sought to combine two of the most well-known experimental philosophy paradigms, the side-effect effect and free will, and examine their joint effects and possible interactions. In doing so, we aimed to extend our understanding of both the SEE and free will attributions in several ways.

First, we tested for the SEE on the ratings of the attribution of free will, asking: Do people attribute higher free will to harmful side-effects than to beneficial side-effects of an action? Second, we tested associations between free will attributions and attributions of responsibility (both blame and praise). Third, we investigated whether manipulating free will universe impacts the SEE. Fourth, we examined whether the two manipulations of (in)determinism and valence are additive or interact to impact attributions of intent, free will, knowledge, regret, and moral responsibility.

Finally, we extended the typical SEE procedure. In the classic SEE paradigm participants read a scenario in which the protagonist is either blamed for harmful actions or praised for helpful actions. To strengthen the between-subject design to also include a within-subject design, within each outcome condition we had participants respond to both praise and blame for both the positive and the negative side-effects. For example, participants who read the SEE scenario that led to a helpful outcome rated both praise for the described positive outcome and blame in case the outcome was different and led to harm. We summarized our hypotheses in [Table 1](#).

Table 1. Summary of hypotheses, rationale, and findings in Studies 1 and 2

Context	H#	Hypothesis	Rationale	Type	Study 1	Study 2
Side-effect effect	1a	Blame attributions for harm > praise attributions for help	Classic side-effect effect.	Confirmatory replication	Supported $d = 1.39$ [1.14, 1.64]	Supported $d = 1.50$ [1.36, 1.63]
	1b	Intent attributions for harm > Intent attributions for help	Classic side-effect effect Blame requires intentionality and causality (Malle et al., 2014), which is not the case for praise (Anderson et al., 2020), thus we can expect the same pattern.	Confirmatory replication	Supported $d = 1.52$ [1.27, 1.77]	Supported $d = 1.61$ [1.48, 1.75]
	1c	Knowledge attributions for harm > Knowledge attributions for help	Beebe and Jensen (2012) found that knowledge is more attributed for harm than for help.	Confirmatory replication	Supported $d = 0.81$ [0.58, 1.04]	Supported $d = 0.73$ [0.61, 0.86]
	2	Free will attributions for harm > Free will attributions for help	Free will has a positive relationship with blame for harm (Feldman et al., 2016; Fillon et al., 2022; Genschow & Vehlow, 2021), but to our knowledge, no investigation was conducted regarding praise. Also, we can draw a direct link with the bad is freer than good (Feldman et al., 2016) concept.	Exploratory	Unsupported $d = 0.12$ [-0.10, 0.34]	Supported $d = 0.18$ [0.06, 0.30].
Interaction with the Universe	3	The SEE effect on blame/praise, intention, and knowledge is weaker in the deterministic universe than in the indeterministic universe.	Based on the possibility that perceptions of agency, or free will, underlie the SEE.	Exploratory	Supported blame/praise $d = -0.82$ [-1.10, -0.53], intention $d = -0.29$ [-0.58, -0.02] and knowledge $d = -0.32$ [-0.60, -0.04]	Supported blame/praise $d = -0.49$ [-0.63, -0.34] intention $d = -0.11$ [-0.25, -0.04], and not supported knowledge $d = -0.01$ [-0.13, 0.16]
SEE at the individual level	4	Blame is more attributed than praise, regardless of the SEE outcome.	Bad is stronger than good (Baumeister et al., 2001). Otherwise, this is the first time, to our knowledge, that blame is assessed for a helpful outcome and praise for a harmful outcome.	Exploratory	Supported $\eta^2_p = 0.41$	Supported $\eta^2_p = 0.38$
Correlations	5a	Free will attributions differ from intent attributions – Free will attributions are weakly or not significantly correlated with intent attributions.	Based on Feldman (2017) Intent and free will are different in nature and are related by the necessity to blame someone.	Exploratory	Supported $r = .15$ [.04, .26]	Supported $r = .08$ [.02, .14]
	5b	Blame attributions are positively correlated to free will attributions.	Based on Malle et al. (2014) and Feldman (2017), free will is a condition to blame and thus, should be positively correlated.	Confirmatory	Supported $r = .54$ [.46, .61]	Supported $r = .50$ [.46, .54]

Context	H#	Hypothesis	Rationale	Type	Study 1	Study 2
	5c	Blame attributions are positively associated with free will attributions, even after controlling for attributions of intent.	Figure 2 from Malle et al. (2014) indicates that intent modulates the relationship between causality and blame, while Table 2 from Feldman (2017) suggests that intentionality is not of the same nature as free will and should not be a necessary condition for the relationship between blame and free will.	Exploratory	Supported <i>r</i> = .33	Supported <i>r</i> = .50 [.45, .54]
Regret (Study 2)	6	Regret attributions for a negative outcome to an agent in the indeterministic universe is higher in comparison to an agent in the deterministic universe.	There is an association between free will and responsibility/blame, we therefore expect that agents in an indeterministic universe will be rated as experiencing higher regret over negative outcomes in comparison to agents in a deterministic universe due to negative side-effect.	Exploratory	N/A	Unsupported <i>d</i> = 0.07 [-0.08, 0.21]

Note. The hypotheses are not clearly stated in the pre-registration of Study 1. We based this table on the hypotheses written in the pre-registration of study 2. Inconsistent findings across Studies 1 and 2 were marked by italics.

Overview, Open Science, Pre-registrations, and Disclosures

We conducted two experiments to test our predictions. Study 1 formed the initial exploratory investigation and was conducted together with another study (we, therefore, consider this to be an exploratory pre-test, see pre-registration of a combined with other research directions <https://osf.io/embrp/>). In Study 2, we pre-registered the specific predictions and ran a dedicated data collection with a larger sample (<https://osf.io/4n5tk/>). All materials, datasets, and analysis scripts are available on the OSF at <https://osf.io/z3g6d/>.

All studies, participants, measures, manipulations, and exclusions conducted for this investigation are reported, and data collection was completed before hypothesis testing. Tests were two-tailed, and α was set at .05.

Study 1: Exploratory Pre-test

Method

Joint Data Collection with Another Project

Our original hypotheses and measures were included as a part of a prior experiment testing another hypothesis by Feldman and Chandrashekar (2018). In Feldman and Chandrashekar's (2018) study, the core experimental manipulations were of a deterministic versus indeterministic universe, focusing on other key measures of interest, and the additional SEE scenarios were added for exploratory purposes (disclosures in their supplementary materials page 2 read: "The data collection included a second part with an experiment regarding the Knobe (2003) side-effect effect. That experiment is unrelated to the research questions in this manuscript and therefore not included or referenced."). Thus, the results presented in this paper are original, going beyond the findings reported in Feldman and Chandrashekar (2018).

Participants

A total of 427 US American participants were recruited from Amazon Mechanical Turk using CloudResearch (Litman et al., 2017). We employed the following CloudResearch options: Duplicate IP Block, and recruited participants with approval rate of 95% and above and who had more than 100 tasks approved. We first excluded 13 participants who indicated a low English proficiency or self-reported not being serious about filling in the survey. These exclusion criteria were not pre-registered for Study 1, yet we applied it to be consistent with the pre-registered criteria

of Study 2. The exclusion criteria did not have much impact and did not change any of the conclusions of the study (differences in effect size were smaller than 0.1), and we provided the results without exclusions with our code. Second, we excluded responses from 101 participants assigned to an additional experimental condition not meant for this investigation¹. Thus, responses from 312 participants were included in this analysis ($M_{age} = 36.2$, $SD_{age} = 12.13$; 179 females). See the supplementary materials for additional details and procedures related to the sample.

Procedure and Design

We summarized the experimental design in Table 2 detailing all the manipulations.

We randomly assigned participants to one of six between-subject conditions in a 3 (universe: deterministic vs. indeterministic universe vs. control) by 2 (negative - harmed the environment vs. positive - helped the environment), first manipulating the presented hypothetical universe and then presenting the classic chairman side-effect effect scenario as taking place in that universe. Manipulations and measures were first pretested in a sample of undergraduates from a university in Hong Kong.

Participants assigned to the deterministic universe and indeterministic universe conditions read a description of the assigned hypothetical universe, then answered comprehension questions and attributions about the described universe to further strengthen the understanding of the described universe. Participants in the universe control condition were not provided with a descriptions of a hypothetical universe. Next, participants were presented with one of the two side-effect effect scenarios. In the deterministic universe and indeterministic universe conditions, the scenarios were described as taking part in the previously described hypothetical universe.

The hypothetical universe related descriptions were adjusted from Nichols and Knobe (2007), which contrasted a fully deterministic universe with a universe in which all is deterministic with the exception of humans. In the original study, the two universes were presented together, yet we adjusted the experimental paradigm to split the two descriptions into two different between-subject conditions. The deterministic and indeterministic universe conditions were presented as follows:

Deterministic universe:

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John

¹ The additional experimental condition presented participants with an "uncertain" universe in which it was not clear whether people are an exception to determinism or not, described and used in Feldman and Chandrashekar (2018). We do not analyze or report this condition as it was not meant for the current investigation, as mirrored by the design of the follow-up Study 2. We also note that because of Feldman and Chandrashekar (2018) the indeterministic and deterministic conditions had several questions more than the control condition that were presented to participants before the side-effect effect scenario and questions.

Table 2. Studies 1 and 2: Experimental design

IV2: Side-effect outcome valence [2 Between]	IV2: Negative outcome: Harmed the environment	IV2: Positive outcome: Helped the environment
IV1: Universe manipulation [3 Between]	[In Universe D there is a company.] The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. <i>Sure enough, the environment was harmed.</i>	[In Universe D there is a company.] The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. <i>Sure enough, the environment was helped.</i>
Control: [No description]	Forced manipulation/comprehension check: "To make sure you understood the scenario - what was the environmental outcome of the chairman's decision to start the new program?" The environment was helped / The environment was harmed / The scenario doesn't say	
Deterministic universe: Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.	Attributions dependent variables: All measures: 1 (<i>Strongly Disagree</i>) to 6 (<i>Strongly Agree</i>).	
Indeterministic universe: Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.	Knowledge "[In Universe D.] the chairman knew the implications of the new program on the environment" "[In Universe D.] the chairman understood the implications of the new program on the environment"	
	Intent "[In Universe D.] the chairman intentions were to have such implications of the new program on the environment?" "[In Universe D.] did the chairman intentionally affect the environment?"	
	Free will "[In Universe D.] the chairman was free to choose not to start the new program" (Reversed) "[In Universe D.] the chairman had to choose what he chose, and could not have chosen to do otherwise"	
	Accountability: Praise attributions "[In Universe D.] the chairman should be applauded for his actions if they led to positive outcomes"	
	Accountability: Blame attributions "[In Universe D.] the chairman should be criticized for his actions if they led to the environment being harmed"	
	Regret/joy [Only in Study 2] "[In Universe D.] the chairman would regret his decision if he learned that his actions led to the environment being harmed."	

Note. Participants in the universe control condition only answered the dependent variables and were not provided with a description of a universe.

decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

Indeterministic universe:

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision

in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

Following the manipulation of the universe, participants read a reminder of the hypothetical universe with the Knobe side-effect managerial scenario as if it was taking part in the hypothetical universe.

The managerial scenario was followed by a manipulation check regarding the outcome of the managerial decision—"what was the environmental outcome of the chairman's decision to start the new program?" (1 = *The environment was helped*; 2 = *The environment was harmed*; 3 = *The scenario does not say*).

Measures

We summarized the experimental design in [Table 2](#) with all the measures used.

Participants evaluated the chairman's behavior on knowledge, intent, free will, and accountability (praise and blame), with all measures on a scale from 1 (*Strongly Disagree*) to 6 (*Strongly Agree*).

Knowledge Attributions

Two items measured the attributions of the extent to which the manager described in the scenario knew about and understood the possible implications of the decision - "In Universe D, the chairman knew the implications of the new program on the environment" and "In Universe D, the chairman understood the implications of the new program on the environment" ($\alpha = .87$).

Intent Attributions

Two items measured attributions of intentionality, whether the manager intended for the program to have the outcome that it did—"In Universe D, the chairman intentions were to have such implications of the new program on the environment?"² and "In Universe D, did the chairman intentionally affect the environment?" ($\alpha = .81$).

Free Will Attributions

Two items measured attributions of free will, whether the manager had the freedom to choose otherwise—"In Universe D, the chairman was free to choose not to start the new program" and "In Universe D, the chairman had to choose what he chose, and could not have chosen to do otherwise" (reversed) ($\alpha = .87$).

Accountability

In the classic SEE experiment, participants typically rated a single dependent variable varied according to the condition, meaning that in the harm condition participants measure measuring blame for the harmful event and praise for the helpful event, we measured blame and praise for both conditions.

Praise attributions

Regardless of the assigned outcome condition, participants rated whether *positive* side-effects deserve *praise* - "In Universe D, the chairman should be applauded for his actions if they led to positive outcomes."

Blame attributions

Regardless of the assigned outcome condition, participants rated whether *negative* side-effects deserve *blame* - "In Universe D, the chairman should be criticized for his actions if they led to the environment being harmed."

Results

Data Analysis

We initially pre-registered an analysis with "Two-way ANOVA with *t*-test contrasts for universe with harm/help" without indicating the type of *t*-test and overlooking the control condition. In our analyses, we used the Welch *t*-test instead of the Student *t*-test because it is more robust to violation of various statistical assumptions (Delacre et al., 2017). We also reported results based on the broader 3 (Experimental condition: Control, Deterministic universe, Indeterministic universe) \times 2 (Outcome: Harm vs. Help) ANOVA. The results based on the pre-registered 2 \times 2 ANOVA are provided in the supplementary materials (with the exclusion of the control condition, as pre-registered). We also ran an exploratory mixed ANOVA with 2 (measures, within: Blame vs. Praise) \times 3 (universe, between: Control vs. Deterministic vs. Indeterministic) \times 2 (outcome, between: Harm vs. Help).

We summarized the descriptive statistics of all dependent variables in six between-subjects experimental conditions in [Table 3](#).

Manipulation Checks

In the harm condition, one participant reported that the chairman helped the environment, and in the help condition, three participants reported that the chairman harmed the environment and three participants selected "the scenario does not say." We also ran the analysis without these participants, and the results were similar without exclusions as for example, the main effect of SEE was $d = 1.59$, 95%CI [1.34, 1.85] for the sample with excluded participants, and $d = 1.52$ [1.27, 1.77] without the exclusions. Because we pre-registered the analysis of the entire sample and the results were similar, below we report results without exclusions.

² A reviewer noted during peer review that the first intent measure was grammatically incorrect. We therefore conducted an analysis for each of the two variables, and found very similar results. Given the appropriate reliability of the two items, we concluded that the participants understood the first sentence as intended, despite the grammatical issues, and proceeded to report the results based on the aggregate. We recommend future research address this issue by rephrasing this specific item.

Table 3. Study 1: Descriptive statistics grouped by experimental conditions

Experimental condition	Outcome	Dimension	Mean	SD
Control	Harm (n = 55)	Praise	2.49	1.60
		Blame	5.36	0.93
		Intention	4.52	1.14
		Free will	5.45	0.91
		Knowledge	5.67	0.55
	Help (n = 56)	Praise	3.07	1.52
		Blame	4.63	1.21
		Intention	2.40	1.16
		Free will	5.26	0.77
		Knowledge	4.88	0.89
Deterministic Universe	Harm (n = 55)	Praise	2.85	1.41
		Blame	3.96	1.57
		Intention	3.95	1.38
		Free will	2.84	1.71
		Knowledge	5.18	1.12
	Help (n = 51)	Praise	2.63	1.39
		Blame	3.45	1.35
		Intention	2.40	1.22
		Free will	2.49	1.49
		Knowledge	4.33	1.43
Indeterministic Universe	Harm (n = 48)	Praise	3.54	1.71
		Blame	5.44	0.77
		Intention	4.54	1.09
		Free will	5.33	1.06
		Knowledge	5.56	0.70
	Help (n = 47)	Praise	3.26	1.28
		Blame	4.36	1.39
		Intention	2.54	1.35
		Free will	5.12	0.84
		Knowledge	4.62	1.25

The Side-effect Effect

Replication of the Original Praise and Blame Effect

We found that blame for a potential negative outcome in the harm condition was higher than praise for a potential positive outcome in the help condition (H1a; $t(307.9) = 12.26$, $p < .001$, $g = 1.39$, 95% CI [1.14, 1.64]). The results were similar for the intentionality (H1b; $t(309.8) = 13.42$, $p < .001$, $g = 1.52$ [1.26, 1.77]) and knowledge (H1c; $t(273.1) = 7.16$, $p < .001$, $g = 0.81$ [0.58, 1.04]). However, we found no support for differences in the free will attributions between the harm and help conditions (H2; $t(309.7) = 1.05$, $p = .294$, $g = 0.12$ [-0.10, 0.34]).

Extension: Differences Between Praise for Positive Outcomes and Blame for Negative Outcomes Regardless of the Assigned Outcome Condition

Participants rated blame for negative outcomes and praise for positive outcomes, regardless of the outcomes in

the scenario. As an exploratory extension, we conducted a 2 (measures, within: Blame vs. Praise) \times 3 (universe, between: Control vs. Deterministic vs. Indeterministic) \times 2 (outcome, between: Harm vs. Help) mixed ANOVA, summarized in Table 5. We found support for an interaction between the measure and the outcome (H4; $F(1, 306) = 13.81$, $p < .001$, $\eta_p^2 = .04$). Praise attributed for a positive side-effect in the harmful outcome condition was similar to the praise attributed for a positive side-effect in the helpful outcome (respectively $M = 2.98$, $SD = 1.42$, $M = 2.94$, $SD = 1.62$) but blame attributed was higher for the negative side-effect in the harmful outcome condition than for a negative side-effect in the helpful outcome condition (respectively $M = 4.90$, $SD = 1.34$, $M = 4.16$, $SD = 1.40$).

Interaction: Indeterminism Manipulation and the Side-effect Effect

The SEE was found across all universe conditions (Figure 1). The effect was the strongest in the indeterministic uni-

verse ($g = 2.05$ [1.53, 2.56]) and the lowest in the deterministic universe (H3; $g = 0.89$ [0.50, 1.29]). The control condition seemed closer to the indeterministic universe condition ($g = 1.80$ [1.35, 2.25]).

The results were similar for the attributions of intention (Figure 2) and knowledge (Figure 3), with stronger effect sizes for the indeterministic and control universes than for the deterministic universe for both intent (indeterministic: $g = 1.16$ [1.15, 2.08]; control: $g = 1.83$ [1.38, 2.27]; deterministic: 1.17 [0.76, 1.58]) and knowledge (indeterministic: $g = 0.92$ [0.49, 1.35]; control: $g = 1.06$ [0.67, 1.47]; deterministic: $g = 0.66$ [0.26, 1.05]). We found no support for a difference between free will attributions (Figure 4) within the universes (indeterministic: $g = 0.23$ [-0.18, 0.62]; control: $g = 0.22$ [-0.15, 0.59]; deterministic: $g = 0.22$ [-0.16, 0.59]).

The results from the ANOVA (Table 4) indicated a main effect of the harmful/helpful outcome manipulation on all variables but free will attribution, a main effect of the type of universe manipulation on all variables but intentionality attribution, and an interaction effect only for the praise/blame attributions.

Extension: Praise and Blame Within-subject Regardless of Assigned Outcome

We tested the interaction between praise and blame at the individual level. We conducted a mixed ANOVA with 2 (measures, within: Blame vs. Praise) \times 3 (universe, between: Control vs. Deterministic vs. Indeterministic) \times 2 (outcome, between: Harm vs. Help). We summarized the results of the ANOVA in Table 5, plotted in Figure 5.

We found support for a main effect of praise and blame, as blame attributions were higher than praise attributions. Praise and blame attributions were higher in the indeterministic universe than in the deterministic, the control group was closer to the deterministic universe for praise, and closer to the indeterministic universe for blame, $F(2, 306) = 11.86$, $p < .001$, $\eta_p^2 = .07$. For praise attributions, we found no support for differences across the universes and the harmful/helpful scenarios. For blame attributions, the harmful and helpful scenarios led to stronger attributions in the indeterministic and control universes than in the deterministic universe. Finally, we found no support for a 3-way interaction ($F(2, 306) = 2.02$, $p = .134$, $\eta_p^2 = 0.002$).

Associations Between Free Will and Blame/praise Attributions

We found support for a positive correlation between free will and blame attributions. For the control group, the correlation was $r(111) = .36$ [.18, .51], $p < .01$, and even stronger when considering the whole sample (H5b), $r(312) = .54$ [.46, .61], $p < .001$. The association held when we ran a partial correlation analysis controlling for the effect of intent and knowledge attributions (H5c; for the control group, $r(111) = .33$, $p < .001$; for the whole sample, $r(312) = .52$, $p < .001$).

We reported the correlations among other attributes in Table 6 (and Table S8 for the correlations by type of universe). Overall, free will attributions had a relatively weak positive correlation with attributions of intent (H5a; $r = .15$

[.04, .26]) and knowledge ($r = .14$ [.03, .24]), and there was a positive correlation between attributions of intent and knowledge ($r = .31$ [.21, .41]). For praise, we only found support for an association with intent ($r = .14$ [.03, .24], no support for a correlation found in the subsample of the control condition).

Discussion

In Study 1, we replicated and extended the well-known findings of the perceived blame/praise asymmetry, intentionality and knowledge of side-effects (Knobe, 2003). Participants attributed more blame for the negative side-effect than praise for the positive side-effect, and more intentionality and knowledge for harm than for help. In line with our predictions, these differences were stronger when the incident was described to be occurring in an indeterministic universe than in a deterministic universe. However, we found no support for differences in free will attributions.

In exploratory extensions, we found that participants attributed more blame than praise for side-effects, regardless of the scenario, and more blame in the indeterministic universe than the deterministic universe, which was not the case for praise. We found that free will attributions were most strongly correlated with blame attributions, even after controlling for ratings of intentionality and knowledge.

Study 2: Confirmatory Investigation

Study 2 was designed to test the robustness of the results noted in Study 1, with a dedicated pre-registration and using a larger well-powered sample. We also added a measure of regret. Based on a priori power analysis, we planned to recruit 1086 participants, with a statistical power of 0.95, an α set to .05, and an effect size of Cohen's $d = 0.20$. The smallest effect size of interest was based on Study 1, which compared free will attribution for a harmful outcome between the indeterministic and deterministic universes.

Method

Participants, Procedures, and Measures

A total of 1108 US American participants were recruited from Amazon Mechanical Turk using CloudResearch (Litman et al., 2017). We employed the following CloudResearch options: Duplicate IP Block, Block Suspicious Geocode Locations, and Verify Worker Country Location, and recruited participants with approval rate of 95% and above and with 1000-500000 approved tasks. After excluding 15 participants following the pre-registered exclusion criteria (see supplementary material for details), the final sample was 1093 (577 females; $M_{\text{age}} = 38.34$, $SD_{\text{age}} = 12.09$).

As in Study 1, we assigned participants randomly to one condition in a 3 (universe manipulation: Indeterministic, Deterministic, Control) \times 2 (outcome: harm vs. help) between-subject design. The scenario descriptions and measures of free will, blame, intentionality, and knowledge at-

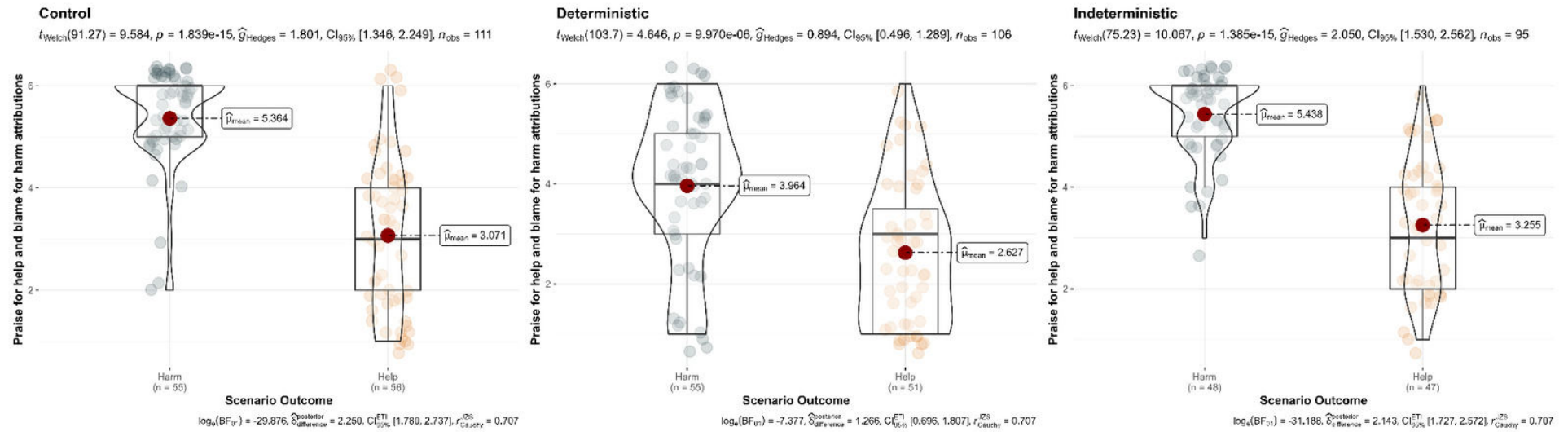


Figure 1. Study 1: Praise/blame attributions across universes (replication of side-effect effect)

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

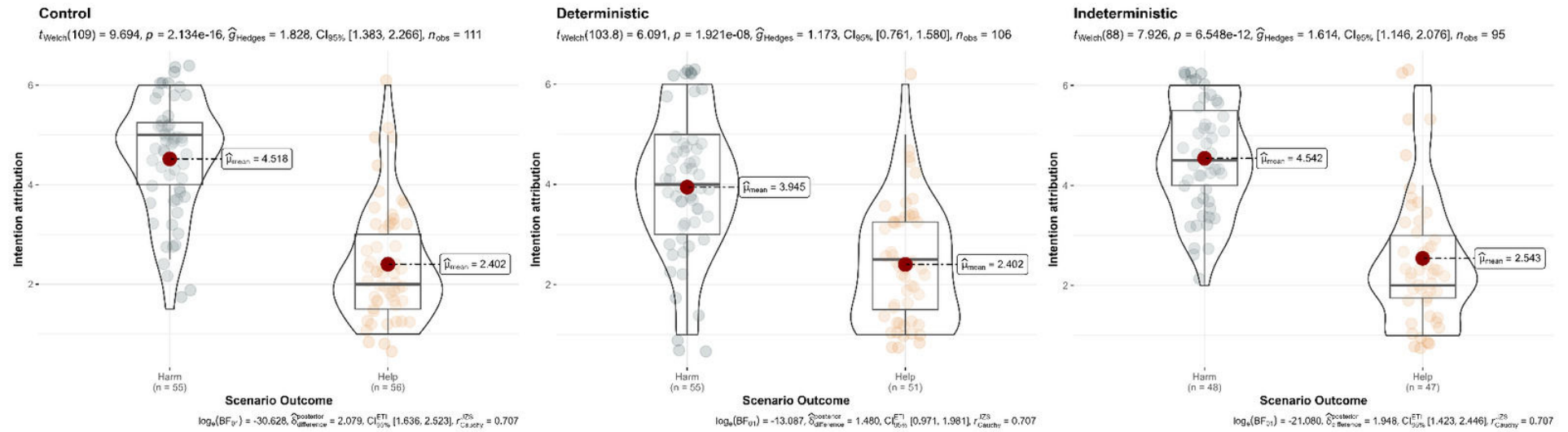


Figure 2. Study 1: Intent attributions across universes

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

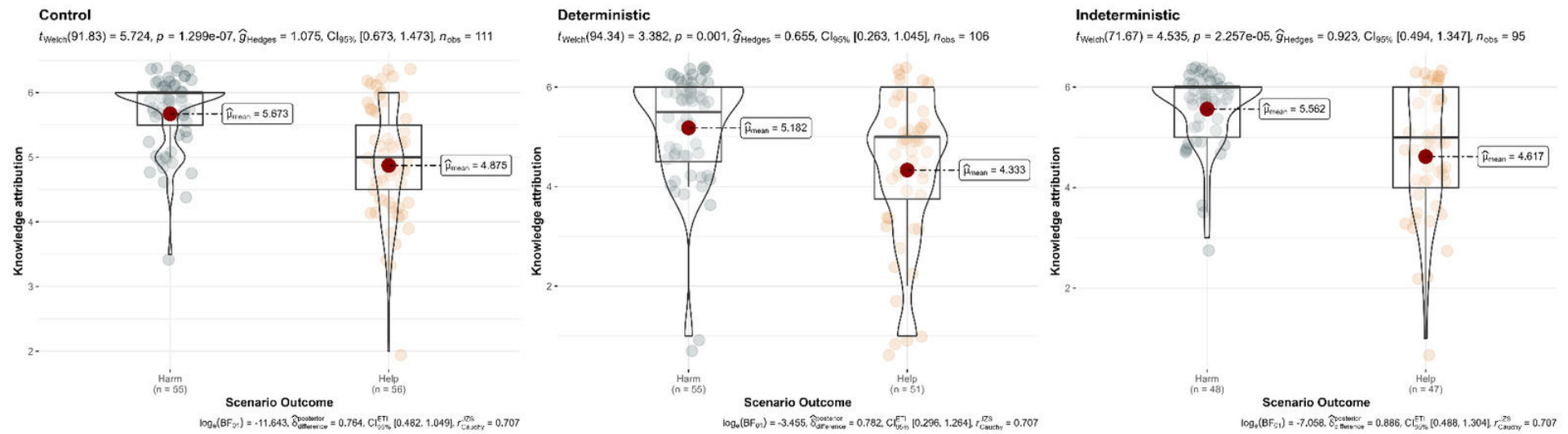


Figure 3. Study 1: Knowledge attributions across universes

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

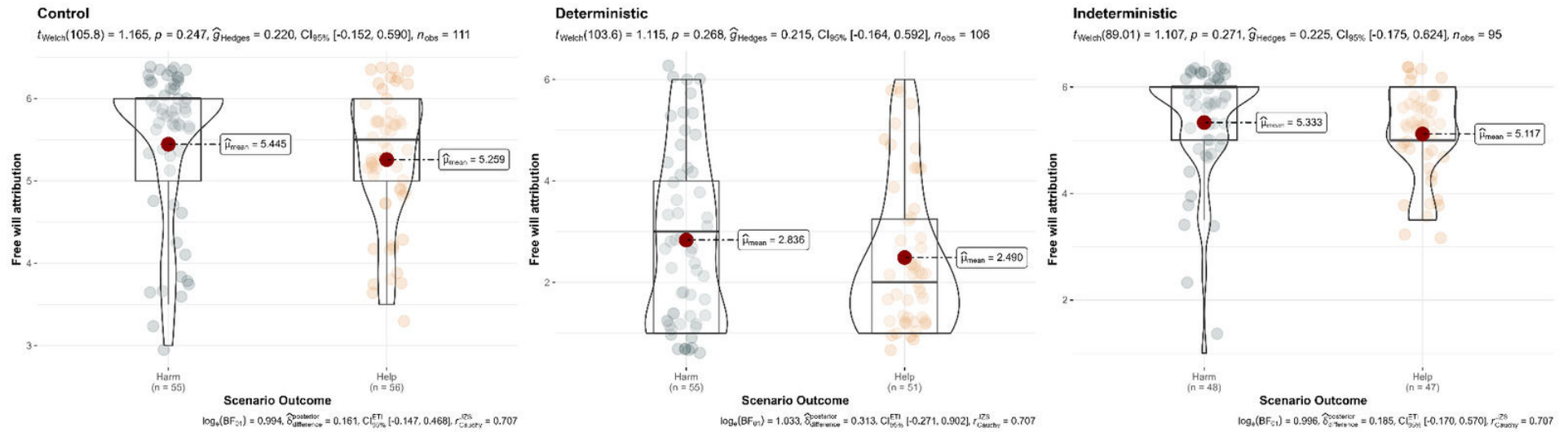


Figure 4. Study 1: Free will attributions across universe conditions

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Table 4. Study 1: Outcome and universe two-way ANOVA for attributions of blame/praise, intentionality, knowledge, and free will

Praise for help and blame for harm attributions						Intentionality attribution					Knowledge attribution					Free Will attribution				
Factor	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p
Outcome (Help vs Harm)	175.75	1	291.21	<.001	.04	182.80	1	276.22	<.001	.37	54.69	1	57.93	<.001	.15	3.44	1	4.84	.065	.01
Universe	20.58	2	34.10	<.001	.12	2.55	2	3.85	.080	.02	6.96	2	7.37	.001	.04	172.27	2	242.29	<.001	.53
Outcome × Universe	4.37	2	7.24	.013	.03	1.62	2	2.44	.201	.01	0.13	2	0.14	.874	.00	0.14	2	0.19	.873	.00

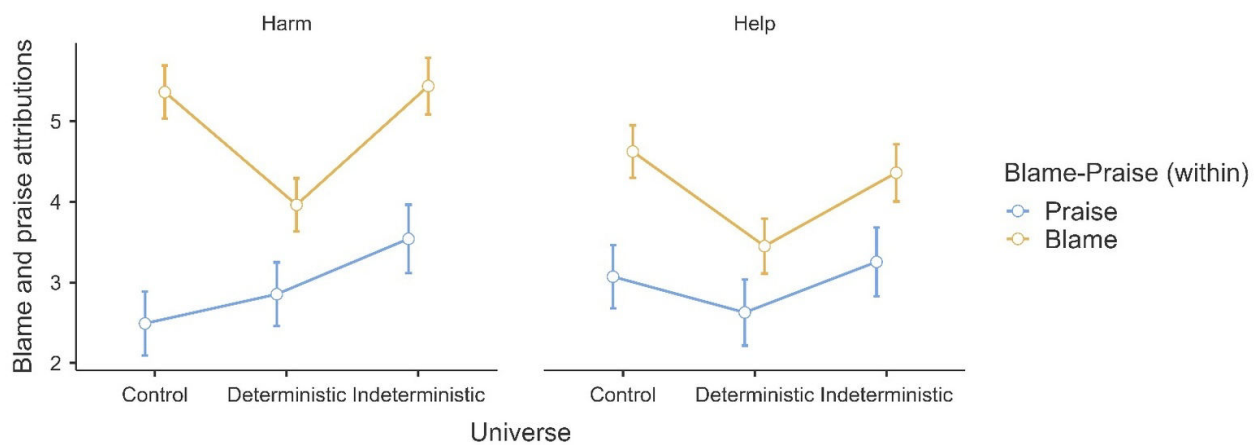
Note. Outcome and Universe are between subject variables. *df* = degree of freedom, *MS* = Mean square, η^2p = partial eta-squared.

Downloaded from http://online.ucpress.edu/collabra/article-pdf/11/1/128423/863738/collabra_2025_11_1_128423.pdf by guest on 06 February 2025

Table 5. Study 1: Praise and blame attributions – 3 way mixed ANOVA testing the effects of measures, outcome, and universe

Factor	Praise/Blame attributions				
	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η_p^2
Measure (Blame vs. Praise)	211.09	1	377.92	<.001	.41
Measure × Outcome (Help vs. Harm)	13.81	1	24.72	<.001	.04
Measure × Universe	11.86	2	21.23	<.001	.07
Measure × Outcome × Universe	2.02	2	3.62	.134	.01

Note. Mixed ANOVA design: 2 (measures, within: Blame vs. Praise) × 3 (universe, between: Control vs. Deterministic vs. Indeterministic) × 2 (outcome, between: Harm vs. Help). *df* = degree of freedom, *MS* = Mean Square, η_p^2 = partial eta-squared.

**Figure 5. Study 1: Estimated marginal means for praise/blame attributions – 3 way mixed ANOVA testing the effects of measures, outcome, and universe****Table 6. Study 1: Means, standard deviations, and correlations across all conditions with confidence intervals**

Variable	<i>M</i>	<i>SD</i>	1	2	3	4	5
1. Free will attributions	4.40	1.72	(.87)				
2. Intent attributions	3.40	1.55	.15** [.04, .26]	(.81)			
3. Knowledge attributions	5.05	1.13	.14* [.03, .24]	.35** [.25, .44]	(.87)		
4. Praise attributions	2.96	1.52	.11 [-.00, .22]	.14* [.03, .24]	-.02 [-.13, .09]	--	
5. Blame attributions	4.53	1.42	.54** [.46, .61]	.32** [.22, .42]	.31** [.21, .41]	.08 [-.04, .19]	--

Note. *M* and *SD* are used to represent mean and standard deviation, respectively. Correlation reported are Spearman correlations. Values in square brackets indicate 95% confidence intervals for each correlation. Alpha coefficients for scales measured with two or more items are on the diagonal cells. *N* = 312. The correlational table by type of universe can be found in Table S9.

tributions were exactly the same as the ones in Study 1, summarized in Table 2. **Added Measure: Attributions of Regret.**

We added a measure of attributions of regret to the agent with one item on a 7-point scale (0 = *Strongly disagree*, 6 = *Strongly agree*). The item was “Do you agree with the following statement? - In Universe D, the chairman would re-

gret his decision if he learned that his actions led to the environment being harmed.”

Results

We summarized descriptive statistics in [Table 7](#). The comprehension check showed that five participants thought the environment was helped in the harm condition, 16 participants thought the environment was harmed in the help condition, and 4 reported that the scenario did not indicate. Exclusions had little to no impact on the findings, for example, the difference between blame/praise was $d = 1.69$ [1.55, 1.83] for the pre-exclusion sample and $d = 1.61$ [1.48, 1.75] post-exclusion. Because we pre-registered the analysis of the entire sample and the results were similar, here we report results without any exclusions.

The Side-effect Effect

Replication of the Original Praise and Blame effect

We found that blame for negative side-effect in the harm condition was higher than praise for a positive side-effect in the help condition (H1a; $t(1088) = 24.75, p < .001, g = 1.50$ [1.36, 1.63]). We found similar results for the intentionality (H1b; $t(1088) = 26.65, p < .001, g = 1.61$ [1.48, 1.75]) and knowledge (H1c; $t(938.4) = 12.11, p < .001, g = 0.73$ [0.61, 0.86]) attributions. Contrary to study 1, we found a smaller difference in free will attributions between the harm and help conditions (H2; $t(1089) = 2.96, p = .003, g = 0.18$ [0.06, 0.30]). Finally, we found no support for a difference concerning our regret extension hypothesis (H6; $t(1091) = 0.34, p = .73, d = 0.02$ [-0.10, 0.14]).

Extension: Differences Between Praise and Blame for Both Helpful and Harmful Side-effects

We measured how participants attributed blame to a negative side-effect of a helpful outcome, and praise to a negative side-effect of a harmful outcome. We conducted a mixed ANOVA with 2 (measures, within: Blame vs. Praise) \times 3 (universe, between: Control vs. Deterministic vs. Indeterministic) \times 2 (outcome, between: Harm vs. Help) reported in [Table 9](#). We found support for an interaction between the measure and the outcome (H4; $F(1, 1087) = 55.8, p < .001, \eta^2_p = .05$). Praise attributed for a positive side-effect in the harmful outcome condition ($M = 2.91, SD = 1.76$) was similar to the praise attributed for a positive side-effect in the helpful outcome condition ($M = 2.99, SD = 1.38$), yet blame attributed to the negative side-effect in the harmful outcome condition was higher ($M = 5.02, SD = 1.32$) than in the helpful outcome condition ($M = 4.20, SD = 1.53$).

Interaction: Indeterminism Manipulation and the Side-effect Effect

As in Study 1, the difference between blame and praise for side-effects was the highest in the indeterministic universe (H3; $g = 2.24$ [1.96, 2.51]), control in between ($g = 2.04$ [1.78, 2.30]) and the lowest in the deterministic universe ($g = 0.78$ [0.57, 1.00]), summarized in [Figure 6](#).

We found similar results for the replication of the SEE, on the attributions of intention ([Figure 7](#)) and knowledge ([Figure 8](#)), with stronger effect sizes for the indeterministic and control universes than for the deterministic universe for both intent (indeterministic: $g = 1.81$ [1.57, 2.06]; control: $g = 1.67$ [1.43, 1.91]; deterministic: $g = 1.37$ [1.14, 1.60]) and knowledge (indeterministic: $g = 0.79$ [0.57, 1.00]; control: $g = 0.74$ [0.53, 0.96]; deterministic: $g = 0.66$ [0.45, 0.87]). There was support for differences between free will attributions ([Figure 9](#)) in the indeterministic and even stronger in the control universes, but no support in the deterministic universe (indeterministic: $g = 0.24$ [0.04, 0.45]; control: $g = 0.68$ [0.47, 0.89]; deterministic: $g = 0.12$ [-0.09, 0.32]). Finally, we found no support for an effect on regret attribution (all $g < 0.16$).

The results from the ANOVA ([Table 8](#); [Figure 10](#)) indicated a main effect of the harmful/helpful scenarios on all variables, a main effect of the type of universe on praise/blame and free will attributions, and an interaction effect only for the praise/blame attributions (excepted regret attributions which were not affected by the scenario and universe for all conditions).

Extension: Praise and Blame Within-subject Regardless of Assigned Outcome

We tested the interaction between praise and blame at the individual level. We conducted a mixed ANOVA with 2 (measures, within: Blame vs. Praise) \times 3 (universe, between: Control vs. Deterministic vs. Indeterministic) \times 2 (outcome, between: Harm vs. Help), and summarized the results in [Table 9](#), plotted in [Figure 11](#).

We found a main effect of praise-blame, as blame was attributed with more intensity than praise for side-effects of both outcomes, $F(1, 1087) = 744.1, MS = 1486.60, p < .001$. Praise was as likely attributed in the indeterministic universe than in the deterministic universe, the control group having a lower praise attribution than the other universes. On the other side, blame was higher in the indeterministic and control universes than in the deterministic universe. For the praise attribution, we found no support for a difference across the universes and the harmful/helpful scenarios. For the blame attribution, the harmful and helpful scenarios led to a stronger attribution in the indeterministic and control than in the deterministic universe. Finally, we found support for a 3-way interaction $F(2, 1087) = 17.2, p < .001, \eta^2_p = 0.03$.

Relationships Between Free Will and Blame Attributions

We tested the associations between free will and blame attributions. As we hypothesized, we found strong support for a positive correlation between free will attribution and blame ratings: for the control group: $r(365) = .48$ [.40, .56], $p < .001$; for the whole sample (H5b), $r(1091) = .50$ [.46, .54], $p < .001$. The results held after controlling for the intent and knowledge attributions (for the control group: partial $r(363) = .44$ [.36, .52]; $p < .001$; for the whole sample (H5c), partial $r(1089) = .50$ [.45, .54]; $p < .001$). We reported

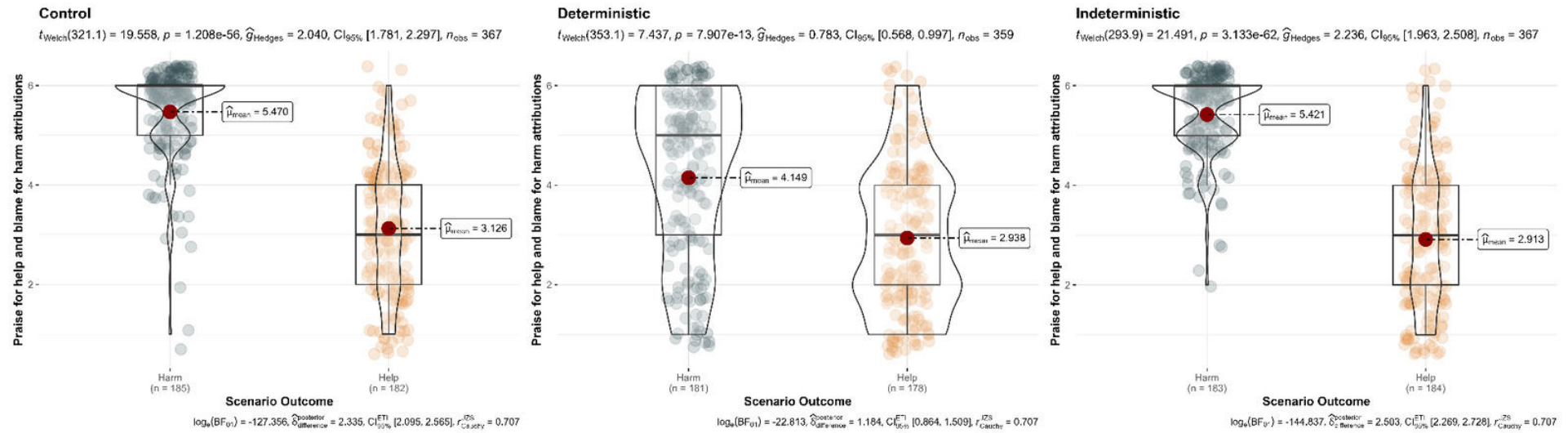


Figure 6. Study 2: Praise/blame attributions across universes (replication of side-effect effect)

Note. A 3 (between subject; universe: control vs. deterministic vs. indeterministic) by 2 (between subject; outcome: harm vs. help) violin plots of praise/blame attributions. To mirror the classic side-effect effects, in this figure the dependent variable varies between the conditions, with blame attributions for the harm condition and praise attributions for the help condition.

Boxplots display the median, first, and third quartiles, and the red circle indicated the mean value.

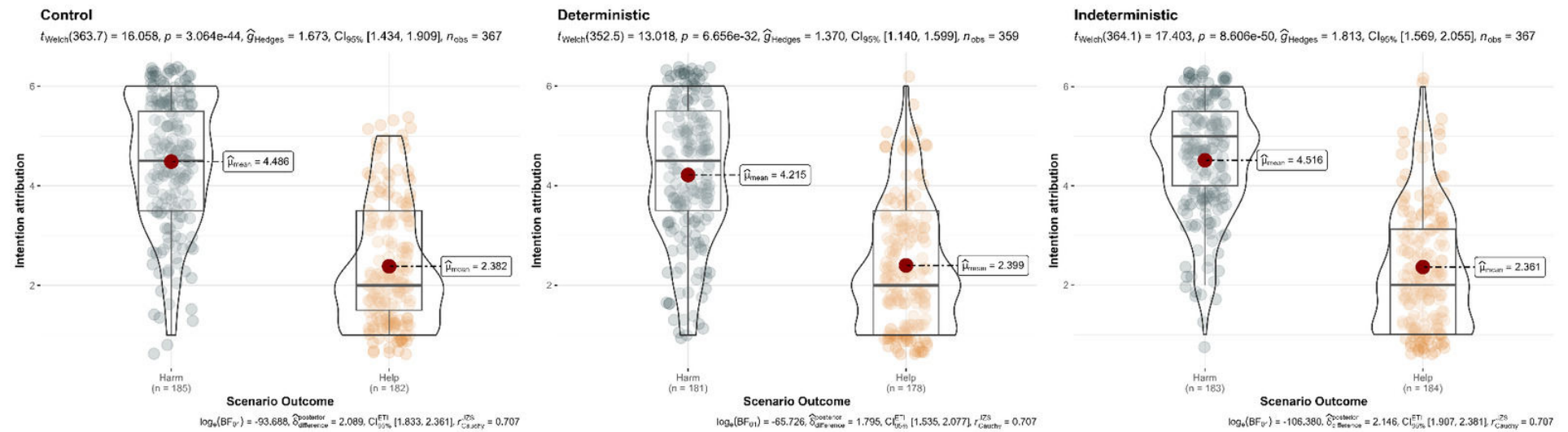


Figure 7. Study 2: Intentionality attributions across universe conditions

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

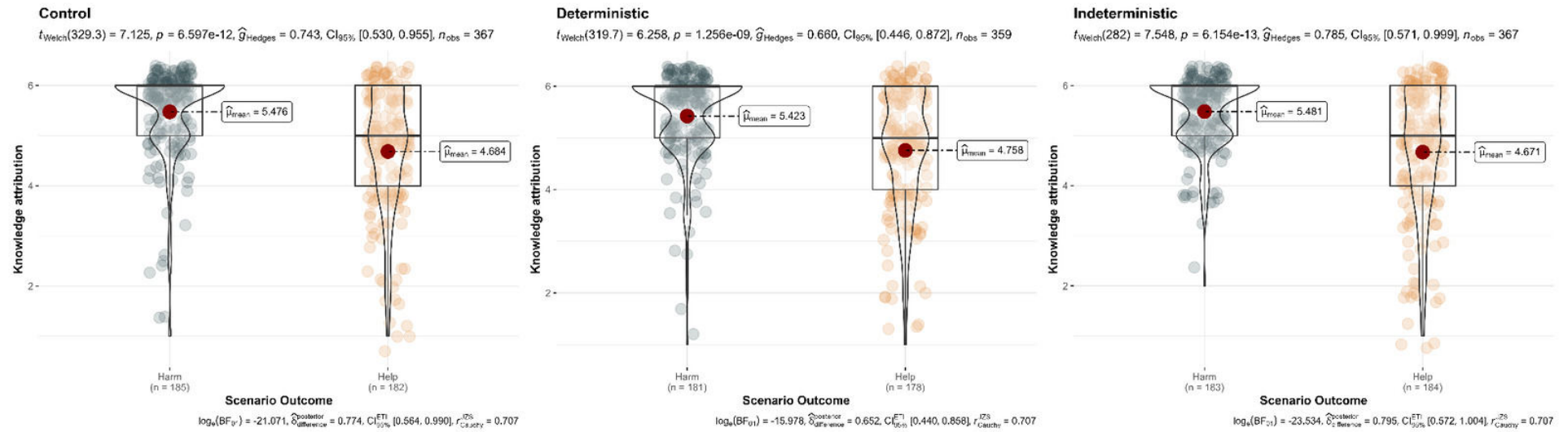


Figure 8. Study 2: Knowledge attributions across universe conditions

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

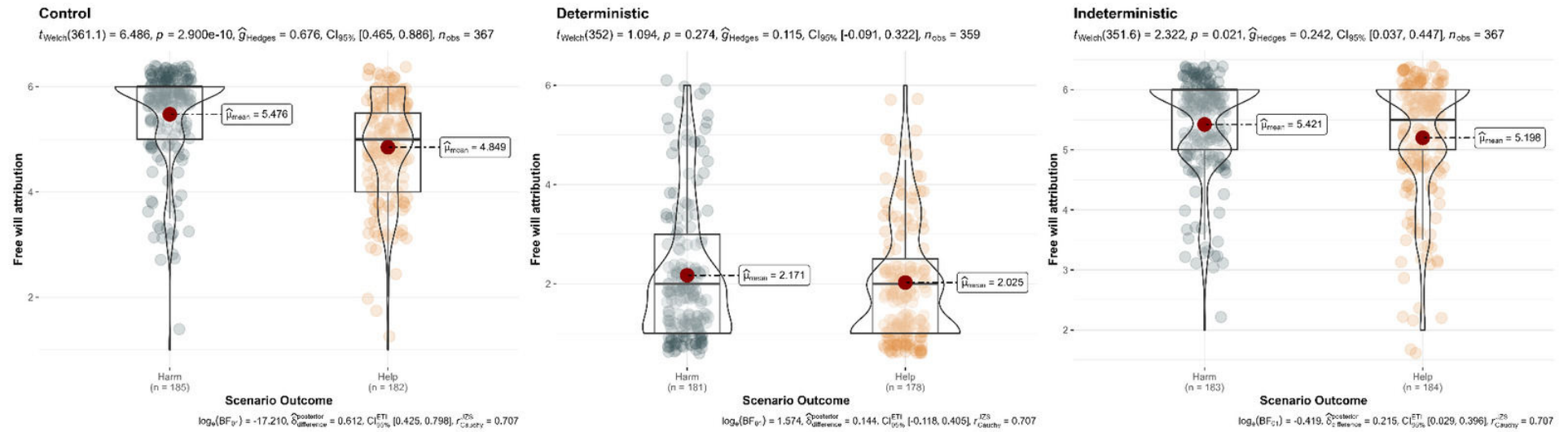


Figure 9. Study 2: Free will attributions across universe conditions

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Table 7. Study 2: Descriptive statistics grouped by experimental conditions

Experimental condition	Outcome	Dimension	Mean	SD
Control	Harm (n = 185)	Praise	2.16	1.54
		Blame	5.47	0.92
		Intention	4.49	1.30
		Free will	5.48	0.88
		Knowledge	5.48	0.88
		Regret	2.09	1.39
	Help (n = 182)	Praise	3.13	1.33
		Blame	4.54	1.26
		Intention	2.38	1.21
		Free will	4.85	0.97
		Knowledge	4.68	1.22
		Regret	2.23	1.29
Deterministic Universe	Harm (n = 181)	Praise	3.08	1.63
		Blame	4.15	1.63
		Intention	4.22	1.41
		Free will	2.17	1.35
		Knowledge	5.42	0.82
		Regret	2.40	1.44
	Help (n = 178)	Praise	2.94	1.45
		Blame	3.49	1.60
		Intention	2.40	1.23
		Free will	2.03	1.18
		Knowledge	4.76	1.16
		Regret	2.17	1.27
Indeterministic Universe	Harm (n = 183)	Praise	3.50	1.83
		Blame	5.42	0.79
		Intention	4.52	1.15
		Free will	5.42	0.82
		Knowledge	5.48	0.69
		Regret	2.20	1.32
	Help (n = 184)	Praise	2.91	1.37
		Blame	4.53	1.49
		Intention	2.36	1.22
		Free will	5.20	1.01
		Knowledge	4.67	1.28
		Regret	2.20	1.23

the correlations among other attributes in [Table 10](#) for the whole sample, and in Table S11 for the correlations per each of the universe conditions. Overall, free will attributions had a weaker positive correlation with attributions of intent (H5a; $r = .08$ [.02, .14]; $p < .01$), knowledge ($r = .10$ [.04, .15]; $p < .01$), and negative with regret ($r = -.07$ [-.13, -.01]; $p = .014$).

Discussion

Study 2 results were largely consistent with the findings of Study 1. In line with our predictions, we found support for the side-effect effect in attributions of praise/blame, ex-

tended to outcome asymmetries regarding intent, knowledge, and to a smaller size, free will attributions. We also found support for blame, intent, knowledge, and free will attributions as being affected by (in)determinism. Finally, we found support for an interaction between (in)determinism and outcome for free will and praise/blame attributions, but not for intent and knowledge. We found no support for regret attributions effects. Finally, we found a higher magnitude of blame attribution than praise for the same participant, who attributed more blame in the indeterministic and control universes than the deterministic

Table 8. Study 2: Outcome and universe two-way ANOVA for attributions of blame/praise, intentionality, knowledge, free will, and regret

Praise/Blame						Intentionality					Knowledge					Free will					Regret				
Factor	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p
Outcome (Help vs Harm)	678.8	1	1115.75	< .001	.38	710.41	1	1120.76	<.001	.40	146.41	1	155.81	<.001	.12	27.43	1	30.06	<.001	0.025	0.13	1	0.22	.72	0.00
Universe	35.8	2	58.79	<.001	.06	1.28	2	4.04	.28	.002	0.19	2	0.02	.98	.001	1084.69	2	1188.65	<.001	0.667	0.87	2	1.54	.42	0.002
Outcome x Universe	27.4	2	45.09	< .001	.05	1.91	2	3.02	.148	.004	0.53	2	0.57	.586	.001	5.56	2	6.09	<.001	0.01	1.68	2	2.96	.19	0.003

Note. Outcome and universe are both between-subject manipulations. *df* = degree of freedom, *MS* = Mean Square, η^2p = partial eta-squared.

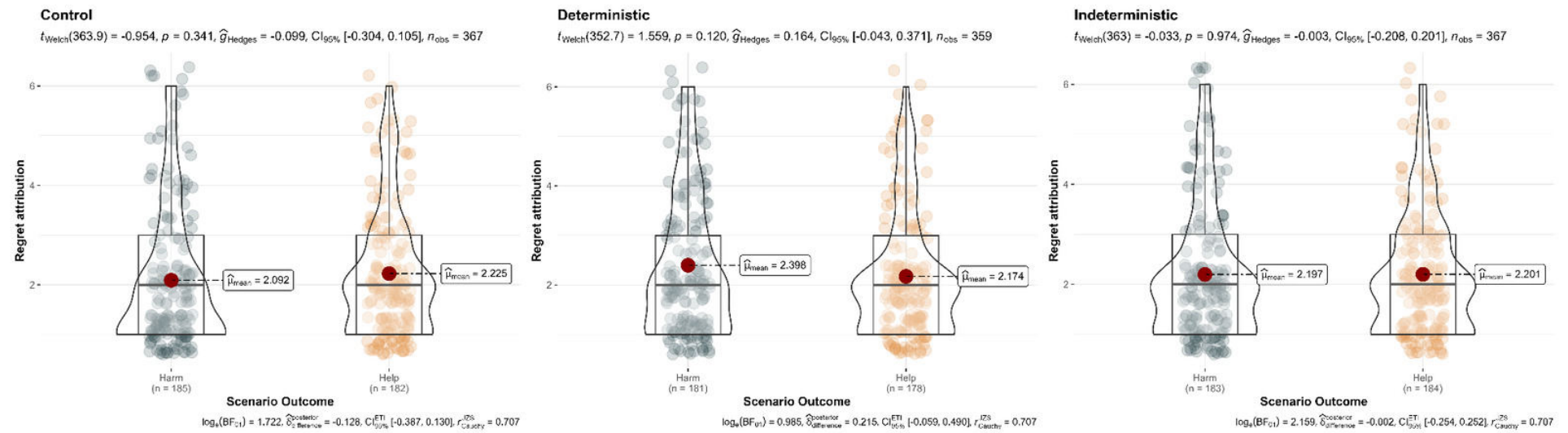


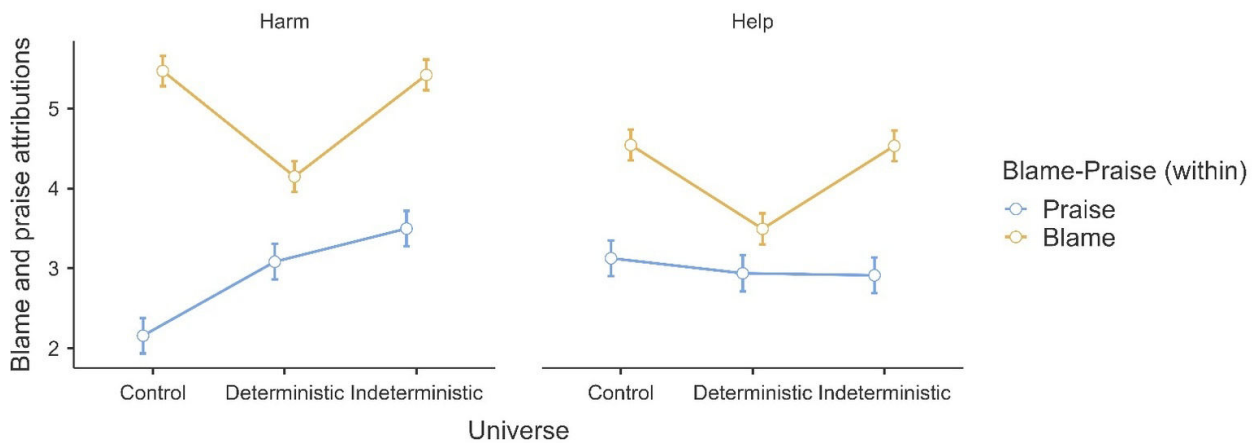
Figure 10. Study 2: Regret attributions across universe conditions

Note. Violin plots of the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Table 9. Study 2: Praise and blame attributions - Results of mixed ANOVA testing the effects of measures, outcome, and universe

Factor	Praise/Blame attributions				
	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2_p
Measure (Blame vs Praise)	744.1	1	1486.60	<.001	.13
Measure × Outcome (Help vs Harm)	55.8	1	111.47	<.001	.10
Measure × Universe	55.8	2	111.46	<.001	.20
Measure × Outcome × Universe	17.2	2	34.31	<.001	0.006

Note. Measures is a within-subject manipulation, outcome and universe are between-subject manipulation. *df* = degree of freedom, *MS* = Mean Square, η^2_p = partial eta squared.

**Figure 11. Study 2: Estimated marginal means for praise/blame attributions – 3 way mixed ANOVA testing the effects of measures, outcome, and universe**

Note. A 3 (between subject; universe: control vs. deterministic vs. indeterministic) by 2 (between subject; outcome: harm vs. help) by 2 (within subject; outcome: harm vs. help) plots of attributions. In this plot, praise and blame attributions are displayed separately.

Table 10. Study 2 means, standard deviations, and correlations with confidence intervals

Variable	<i>M</i>	<i>SD</i>	1	2	3	4	5
1. Free will attribution	4.21	1.82	(.91)				
2. Intention attribution	3.40	1.61	.08** [.02, .14]	(.84)			
3. Knowledge attribution	5.08	1.10	.10** [.04, .15]	.29** [.23, .34]	(.90)		
4. Praise attribution	2.95	1.58	-.06 [-.12, .00]	.07* [.02, .13]	-.05 [-.11, .01]	--	
5. Blame attribution	4.61	1.49	.50** [.46, .54]	.27** [.22, .33]	.17** [.11, .23]	-.00 [-.06, .06]	--
6. Regret attribution	2.21	1.33	-.07* [-.13, -.01]	.16** [.11, .22]	.12** [.06, .18]	-.20** [-.26, -.14]	-.07* [-.12, -.01]

Note. *M* and *SD* are used to represent mean and standard deviation, respectively. Values in square brackets indicate the 95% confidence interval for each correlation. Alpha coefficients for scales measured with two or more items are on the diagonal cells. The correlational table by type of universe can be found in Table S11.

universe for both scenarios, and especially for harm, which was not the case for praise.

We note three takeaways from our findings. First, we found that side-effect effects are relatively consistent across contexts that vary on the possibility of free will. Attributions of blame/praise for side-effects, intent, and knowledge had a consistent and larger variation for harmful

outcomes than for helpful outcomes. We also found main effect differences between blame and praise with blame attributions generally higher than praise attributions for side-effects, though this could be a psychometric artifact and so more work is needed to identify the nature of difference. Finally, free will attributions had a stronger link with blame than praise attribution, even after accounting

for the perceived intentionality and knowledge attributed to the described agent.

General Discussion

In two studies, we revisited and combined two classic paradigms in experimental philosophy – the side-effect effect and the impact of free will on moral accountability. We successfully replicated these classic effects, and further extended them by examining their joint effects and interactions, with several important insights. We summarized the results of the investigations in [Table 1](#). Below, we will describe and interpret our results concerning the replication of the two main theories tested, the side-effect effect and the free will relationship with blame, before discussing the new findings linking the two theories in light of our extension regarding ascribing blame for negative side-effect of harmful outcome and praise for positive side-effects of helpful outcome.

Side-effect Effect: Replication

Revisiting the SEE and examining the impact of outcome manipulation, we found that participants attributed higher intent, knowledge, and blame to harmful negative side-effects than helpful positive side-effects of an action. The results were consistent across the two samples.

In addition to the replication of the SEE on praise/blame, we replicated the asymmetry in the attributions of intent and knowledge. Our findings are in line with the extant literature on moral judgments (see Malle, 2021, for a review). For example, a recent theoretical assertion defined intentionality as follows: “people consider that an agent did X *intentionally* to the extent that X was causally dependent on how much the agent wanted X to happen (or not to happen)” (Quillien & German, 2021, p. 1). We also replicated the findings from Beebe and Jensen (2012) that knowledge is also more attributed to a harmful than to a helpful side-effect, indicating that laypeople judge the knowledge of others (here, a chairman) based on the outcome of a decision, as they do for intention.

Free Will Manipulation and Attributions

Type of Universe and Moral Accountability: Replication

We tested whether the classic experimental philosophy of manipulating determinism in a described universe impacts evaluations of moral accountability. We found that the manipulation had a strong impact on the attributions made toward the chairman, regardless of whether the side-effect was positive or negative. Across both Studies 1 and 2, participants 1) made stronger attributions in the indeterministic universe than in the deterministic universe, 2) participants attributed more blame and praise in the indeterministic universe than the deterministic universe, and 3) the attributions of intention and knowledge were stronger in the indeterministic universe, but to a lower extent (the

two lower bound for the effect sizes of intention attributions are close to zero, but significant).

Associations Between Free Will and Accountability Attributions: Replication

Why and when do people blame others? We provided one possible answer by showing a link between blame and free will. In negative outcome situations, blame (Monroe et al., 2014) and free will (Feldman et al., 2016) are both about the capacity to choose, suggesting that blame is due to perceiving harm as a choice. Our results supported this view, as we found a positive and strong correlation between blame and free will attributions. We also found stronger blame attributions in the indeterministic universe than in the deterministic universe, which was the measure that displayed the highest variation between the universes in Study 2. The relationship held when being controlled for the other variables. Our correlational and experimental results are in line with some of the recent work on blame. For example, Genschow and Vehlow (2021) found that free will perception was related to not only the blame toward an offender but also toward the victim, meaning that the blame/free will relationship goes beyond the need for compensation. Put together, these findings indicated that the possibility of having free will to act is related to the blame we attribute to someone, but not the attribution of praise or regret.

The literature has remained unclear regarding how moral judgments are related to our lay assumptions regarding the universe. The current results show that the asymmetry between blame for the side-effects of harm/praise for the side-effects of help is stronger in the indeterministic world than in the deterministic world, but also that the “control” universe, in which nothing is said regarding the law of the universe the chairman is in has the same properties as the indeterministic universe. The results seem to indicate that laypersons tend to view the universe we are in as similar to the indeterministic universe described in the vignette, and that individuals seem to perceive indeterminism by default, or at least “laypersons viewed the universe as allowing for human indeterminism” (Feldman & Chandrashekar, 2018, p. 539). We also found that the relationship between free will and intent is weak in both studies, which supports the idea that free will and intent are separate constructs, that free will is not a prerequisite for intention, and that attributions of free will and intent can lead to blame from a different path (Feldman, 2017).

Associations Between Free Will and the Side-effect Effect

We examined the impact of valence of the outcomes in the classic SEE chairman scenario over free will attributions and found support for free will attribution asymmetries in Study 2, but less so in Study 1. In Study 2, free will attributions were higher for the harmful outcome than for the helpful outcome, though the effect was weaker than the effects of intention, knowledge, praise, and blame. However, we found no support for the asymmetry in Study 1. The in-

consistent results across the two samples may be attributed to the smaller sample in Study 1.

Although a link might be made between the SEE and free will, the effect sizes are much smaller for the free will dimensions (attribution of free will and manipulation of the universes) than for the attributions of blame, intent, knowledge, or scenario manipulation. A very recent debate led to the conclusion that free will is attributed more on the basis of norm-violation accounts than because of intrinsic motivation (Monroe & Ysidron, 2021). In this article, researchers found that participants attributed the same amount of free will to a praiseworthy and blameworthy action. Another recent work has indicated that ignorance is a key factor in explaining attributions for an action (Kirfel & Phillips, 2023). People attributed more intentionality and free will to a norm-violating action when the agent was aware of the consequences of his act. This lack of awareness was left ambiguous in our scenarios, as we simply stated that the chairman “did not care” about the consequences) which might have led participants to not infer free will and intentionality as much as we intended, in the case the chairman was not aware of his actions. Furthermore, in both scenarios, the chairman could have been seen as acting in line with norms ascribed to chairmen (maximizing profits) for some participants, who would not attribute free will related to norm violation. On the other side, some participants might have seen a violation of norms related to the environmental protection (the chairman does not care for the environment), attributing more free will for violating this norm. These two arguments can explain the weaker attribution of free will for the difference between harm and blame than the other attributions. To further understand how free will attribution can vary based on a harmful or helpful outcome, researchers should consider manipulating the salience of the norm-violation related to a positive or negative non-intended outcome, but also the degree of awareness of the consequences of the actor’s behavior.

Free Will and Regret Associations

We found a strong association between free will and blame, and therefore expected that agents in an indeterministic universe would be attributed a stronger experience of regret over negative outcomes compared to agents in a deterministic universe. However, our results failed to find experimental support for this view. This result is surprising, as Fillon et al. (2022) found that regret related to free will across the universe conditions for the exceptionality bias. Still, in our study, the manipulation of the universe did not lead to a change in regret attributions. We did not find support for an association between regret and free will attributions, and we found no support for differences in regret across types of universes.

It is possible that our (in)determinism universe manipulations were not effective enough to influence attributions of regret. Alternatively, the descriptive part of the scenario related to the harmful or helpful outcomes might have driven the entire effect, where participants overlooked the universes when ascribing regret. Theoretically, one can only regret a decision if one can think about a better al-

ternative. Thus, regret is only possible if there are alternative choices; in other words, if there is free will—a view supported by the work of Fillon and colleagues (2022). They manipulated the universes to find an interaction between the exceptionality effect and determinism on regret and found that only exceptionality affects regret. The authors indicated that it could be hard for participants to understand the deterministic universe, as laypeople believe that they have free will even in a deterministic universe, a view held by the majority of people called natural compatibilism (Nadelhoffer et al., 2020). Thus, it is possible that the manipulation of the universe might not be a good operationalization for choice representation, because it is hard for participants to represent the differences between the universes and their consequences. More work is needed to explore these directions.

Blame for Side-Effects is Stronger than Praise Regardless of the Outcome

Related to moral accountability, we found that attributions of blame for potential negative side-effects were stronger than attributions of praise for potential positive side-effects, regardless of the outcomes (harmful or helpful) described in the scenario. Based on the view that bad is stronger than good (Baumeister et al., 2001), we expected people to attribute stronger blame than praise for a potential side-effect. Our findings supported the prediction across both Study 1 ($d = 1.39$) and Study 2 ($d = 1.50$). Interestingly, Feldman et al. (2016) argued that “bad is freer than good.” People attributed higher free will to negative than positive valence, regardless of morality or intent, for both self and others. We found similar though weaker results for side-effects. This finding strengthens the argument for a relationship between blame and free will attributions.

Limitations

One limitation lies in our manipulation of the free will universe, as it refers very broadly to the ability to choose without disentangling the constraints underlying the inability to choose. In the discussion regarding free will, there are context-dependent constraints (e.g., job role) and broader, more fundamental factors that restrict choice that are close to the philosophical debate on free will (e.g., determinism, higher power, genes, physics, etc.). While the free will universe manipulation is close to the philosophical debate and the manipulation impacted free will attributions, it is possible that free will attributions might also be related to the contextual aspects of choice. Finally, there is the possibility that the universe scenarios do not work as intended, as participants can have difficulties understanding the consequences of a deterministic universe. Future studies can expand on our findings to examine more specific constraints and how the effects we reported vary depending on the type of constraint or operationalization of the universe.

We also found an oversight in the two items measuring attribution of intention, which were not grammatically

clear. We adapted the two items on the intention from Jamison and colleagues' (2020) study, which are not standardized and may have impacted one of the questions. However, the reliability coefficients were high for both studies. We conducted a 2x2 ANOVA on both items and noted that, even if the second item is higher than the first, they are affected the same way by the SEE and the type of universe. We added the results to OSF. We believe that despite the awkward phrasing, the two questions were similarly understood by the participants.

Regarding regret, we measured regret attribution as the possibility of experiencing regret for a decision if it led to the environment being harmed. Therefore, this measure was used for helpful and harmful outcomes, and assessed counterfactuals. To answer this question, participants have to think about the action, the possibility for this action to lead to a harmful situation, and then how the chairman should feel regarding these non-existent consequences. Understanding this process is not trivial and requires imagination to construct an answer, especially if the decision leads to helping the environment. Using descriptions might not be the best way to test the relationship between side-effects and regret attributions.

We tested if the chairman's knowledge about the program were associated with blame or praise for side-effects. We asked participants to state how well the chairman knew and understood the implications of his program on the environment. However, knowledge is only a proxy for attribution of responsibility – and causality. Stated differently, if the chairman knew about the side-effects of his acts, it does not automatically mean that he was responsible, or that he wanted to cause these side-effects. A new line of research regarding the relationship between knowledge and causality could draw on our results and ask participants if, by knowing about the side-effects and still choosing to implement the program, the chairman could be seen as responsible, and a cause for the side-effects of his harmful or helpful program.

Conclusion

In two experiments, we combined together two influential theories in experimental philosophy regarding blame: the side-effect effect and the relationship between blame and free will. We successfully replicated the side-effect effect, but also found support for the relationship between blame ascribed to side-effects and free will, with correlational evidence in both studies and the impact of a determinism manipulation of the description of the universe in Study 2. We then found that the relationship between blame and free will was stronger than attributions of intent or knowledge, suggesting that participants blame more freer actions not solely because these actions were intended. Finally, we tested for the first time if the side-effect effect could be seen regardless of the harmful and helpful outcome and found that blame was always attributed more than praise for potential side-effects. This finding is in line with the “bad is stronger than good” but also “bad is freer

than good” hypotheses. Further work is needed to understand the causal relationships between the freeness to act, the attribution of accountability, and blame for unintended consequences of actions.

Competing Interests

The author(s) declared no potential conflicts of interests with respect to the authorship and/or publication of this article.

Funding

Nothing to disclose.

Authorship Declaration

Gilad Feldman designed Study 1. Prasad designed and wrote the pre-registration for follow-up Study 2. Gilad conducted data collection for both studies. Prasad and Adrien analyzed the data. Adrien and Prasad wrote the manuscript and edited revisions. Prasad, Adrien, and Gilad jointly finalized the manuscript for submission.

CRediT (Contributor Roles Taxonomy)

Role	Adrien Fillon	Subramanya Prasad Chandrashekar	Gilad Feldman
Conceptualization		X	X
Pre-registration		X	X
Data curation			X
Formal analysis	X	X	
Funding acquisition			X
Investigation	X	X	
Pre-registration peer review / verification			X
Data analysis peer review/ verification	X	X	
Methodology	X	X	X
Project administration			X
Resources			X
Software			
Supervision			X
Validation	X	X	
Visualization	X	X	
Writing-original draft	X	X	
Writing-review and editing	X	X	X

Submitted: February 05, 2023 PST. Accepted: September 24, 2024 PST. Published: February 04, 2025 PST.



This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CCBY-4.0). View this license's legal deed at <http://creativecommons.org/licenses/by/4.0> and legal code at <http://creativecommons.org/licenses/by/4.0/legalcode> for more information.

References

- Anderson, R. A., Crockett, M. J., & Pizarro, D. A. (2020). A Theory of Moral Praise. *Trends in Cognitive Sciences*, 24(9), 694–703. <https://doi.org/10.1016/j.tics.2020.06.008>
- Baumeister, R. F., Bratslavsky, E., Finkenauer, C., & Vohs, K. D. (2001). Bad is Stronger than Good. *Review of General Psychology*, 5(4), 323370. <https://doi.org/10.1037/1089-2680.5.4.323>
- Beebe, J. R., & Buckwalter, W. (2010). The Epistemic Side-Effect Effect. *Mind & Language*, 25(4), 474498. <https://doi.org/10.1111/j.1468-0017.2010.01398.x>
- Beebe, J. R., & Jensen, M. (2012). Surprising connections between knowledge and action: The robustness of the epistemic side-effect effect. *Philosophical Psychology*, 25(5), 689715. <https://doi.org/10.1080/09515089.2011.622439>
- Chandrashekar, S. P., Chan, Y. Y., Cheng, K. L., Yao, D., Lo, C. Y. S., Cheung, T. C. A., ... Feldman, G. (2022). Revisiting the folk concept of intentionality: Replications of Malle and Knobe (1997). *Journal of Experimental Social Psychology*, 102, 104372. <https://doi.org/10.1016/j.jesp.2022.104372>
- Cokely, E. T., & Feltz, A. (2009). Individual differences, judgment biases, and theory-of-mind: Deconstructing the intentional action side effect asymmetry. *Journal of Research in Personality*, 43(1), 1824. <https://doi.org/10.1016/j.jrp.2008.10.007>
- Cova, F., & Naar, H. (2012). Side-Effect effect without side effects: The pervasive impact of moral considerations on judgments of intentionality. *Philosophical Psychology*, 25(6), 837854. <https://doi.org/10.1080/09515089.2011.622363>
- Cushman, F., Knobe, J., & Sinnott-Armstrong, W. (2008). Moral appraisals affect doing/allowing judgments. *Cognition*, 108(1), 281–289. <https://doi.org/10.1016/j.cognition.2008.02.005>
- Delacre, M., Lakens, D., & Leys, C. (2017). Why Psychologists Should by Default Use Welch's t-test Instead of Student's t-test. *International Review of Social Psychology*, 30(1), 92–101. <https://doi.org/10.5334/irsp.82>
- Feldman, G. (2017). Making sense of agency: Belief in free will as a unique and important construct. *Social and Personality Psychology Compass*, 11(1), e12293. <https://doi.org/10.1111/spc3.12293>
- Feldman, G., & Chandrashekar, S. P. (2018). Laypersons' Beliefs and Intuitions About Free Will and Determinism: New Insights Linking the Social Psychology and Experimental Philosophy Paradigms. *Social Psychological and Personality Science*, 9(5), 539549. <https://doi.org/10.1177/1948550617713254>
- Feldman, G., Wong, K. F. E., & Baumeister, R. F. (2016). Bad is freer than good: Positive–negative asymmetry in attributions of free will. *Consciousness and Cognition*, 42, 2640. <https://doi.org/10.1016/j.concog.2016.03.005>
- Feltz, A. (2007). The Knobe Effect: A Brief Overview. *The Journal of Mind and Behavior*, 28(3/4), 265–277. <http://www.jstor.org/stable/43854197>
- Feltz, A., & Cova, F. (2014). Moral responsibility and free will: A meta-analysis. *Consciousness and Cognition*, 30, 234246. <https://doi.org/10.1016/j.concog.2014.08.012>
- Fillon, A., Lantian, A., Feldman, G., & N'Gbala, A. (2022). Exceptionality Effect in Agency Attributions: Exceptional Behaviors are Perceived as Higher Free will than Routine Behaviors. *International Review of Social Psychology*, 35(1), 3. <https://doi.org/10.5334/irsp.591>
- Genschow, O., & Vehlow, B. (2021). Free to blame? Belief in free will is related to victim blaming. *Consciousness and Cognition*, 88, 103074. <https://doi.org/10.1016/j.concog.2020.103074>
- Guglielmo, S., & Malle, B. F. (2010). Can Unintended Side Effects Be Intentional? Resolving a Controversy Over Intentionality and Morality. *Personality and Social Psychology Bulletin*, 36(12), 1635–1647. <https://doi.org/10.1177/0146167210386733>
- Heider, F. (1958). The naive analysis of action. In *The psychology of interpersonal relations* (pp. 79–124). John Wiley & Sons Inc. <https://doi.org/10.1037/10628-004>
- Hindriks, F. (2008). Intention Action and the Praise-Blame Asymmetry. *The Philosophical Quarterly*, 58(233), 630–641. <https://doi.org/10.1111/j.1467-9213.2007.551.x>
- Jamison, J., Yay, T., & Feldman, G. (2020). Action-inaction asymmetries in moral scenarios: Replication of the omission bias examining morality and blame with extensions linking to causality, intent, and regret. *Journal of Experimental Social Psychology*, 89, Article 103977. <https://doi.org/10.1016/j.jesp.2020.103977>
- Kirfel, L., & Phillips, J. (2023). The pervasive impact of ignorance. *Cognition*, 231, 105316. <https://doi.org/10.1016/j.cognition.2022.105316>
- Klein, R. A., Vianello, M., Hasselman, F., Adams, B. G., Adams, R. B., Alper, S., ... Sowden, W. (2018). Many Labs 2: Investigating variation in replicability across samples and settings. *Advances in Methods and Practices in Psychological Science*, 1(4), 443–490. <https://doi.org/10.1177/2515245918810225>
- Knobe, J. (2003). Intentional action and side effects in ordinary language. *Analysis*, 63(279), 190194. <https://doi.org/10.1111/1467-8284.00419>
- Komatsu, T., Malle, B. F., & Scheutz, M. (2021). Blaming the Reluctant Robot: Parallel Blame Judgments for Robots in Moral Dilemmas across U.S. and Japan. *Proceedings of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, 6372. <https://doi.org/10.1145/3434073.3444672>
- Laurent, S. M., Reich, B. J., & Skorinko, J. L. M. (2019). Reconstructing the side-effect effect: A new way of understanding how moral considerations drive intentionality asymmetries. *Journal of Experimental Psychology: General*, 148(10), 17471766. <https://doi.org/10.1037/xge0000554>

- Litman, L., Robinson, J., & Abberbock, T. (2017). TurkPrime. com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behavior Research Methods*, 49(2), 433–442.
- Malle, B. F. (2021). Moral Judgments. *Annual Review of Psychology*, 72(1), 293318. <https://doi.org/10.1146/annurev-psych-072220-104358>
- Malle, B. F., Guglielmo, S., & Monroe, A. E. (2014). A theory of blame. *Psychological Inquiry*, 25(2), 147–186. <https://doi.org/10.1080/1047840x.2014.877340>
- Malle, B. F., Guglielmo, S., Voiklis, J., & Monroe, A. E. (2022). Cognitive blame is socially shaped. *Current Directions in Psychological Science*, 31(2), 169–176. <https://doi.org/10.1177/09637214211068845>
- Malle, B. F., & Knobe, J. (1997). The folk concept of intentionality. *Journal of Experimental Social Psychology*, 33, 101–121. <https://doi.org/10.1006/jesp.1996.1314>
- Martin, J. W., Jordan, J. J., Rand, D. G., & Cushman, F. (2019). When do we punish people who don't? *Cognition*, 193, 104040. <https://doi.org/10.1016/j.cognition.2019.104040>
- Monroe, A. E., Dillon, K. D., & Malle, B. F. (2014). Bringing free will down to Earth: People's psychological concept of free will and its role in moral judgment. *Consciousness and Cognition*, 27, 100108. <https://doi.org/10.1016/j.concog.2014.04.011>
- Monroe, A. E., & Malle, B. F. (2019). People systematically update moral judgments of blame. *Journal of Personality and Social Psychology*, 116(2), 215236. <https://doi.org/10.1037/pspa0000137>
- Monroe, A. E., & Ysidron, D. W. (2021). Not so motivated after all? Three replication attempts and a theoretical challenge to a morally motivated belief in free will. *Journal of Experimental Psychology: General*, 150(1), e1–e12. <https://doi.org/10.1037/xge0000788>
- Nadelhoffer, T., Rose, D., Buckwalter, W., & Nichols, S. (2020). Natural compatibilism, indeterminism, and intrusive metaphysics. *Cognitive Science*, 44(8), e12873. <https://doi.org/10.1111/cogs.12873>
- Nahmias, E., Coates, D. J., & Kvaran, T. (2007). Free Will, Moral Responsibility, and Mechanism: Experiments on Folk Intuitions. *Midwest Studies In Philosophy*, 31(1), 214242. <https://doi.org/10.1111/j.1475-4975.2007.00158.x>
- Nichols, S., & Knobe, J. (2007). Moral Responsibility and Determinism: The Cognitive Science of Folk Intuitions. *Noûs*, 41(4), 663685. <https://doi.org/10.1111/j.1468-0068.2007.00666.x>
- Ohtsubo, Y. (2007). Perceived intentionality intensifies blameworthiness of negative behaviors: Blame-praise asymmetry in intensification effect 1. *Japanese Psychological Research*, 49(2), 100–110. <https://doi.org/10.1111/j.1468-5884.2007.00337.x>
- Pettit, D., & Knobe, J. (2009). The pervasive impact of moral judgment. *Mind & Language*, 24(5), 586–604. <https://doi.org/10.1111/j.1468-0017.2009.01375.x>
- Phillips, J., & Knobe, J. (2009). Moral judgments and intuitions about freedom. *Psychological Inquiry*, 20(1), 30–36. <https://doi.org/10.1080/10478400902744279>
- Quillien, T., & German, T. C. (2021). A simple definition of 'intentionally'. *Cognition*, 214, 104806. <https://doi.org/10.1016/j.cognition.2021.104806>
- Shultz, T. R., & Wells, D. (1985). Judging the intentionality of action-outcomes. *Developmental Psychology*, 21(1), 83. <https://doi.org/10.1037/0012-1649.21.1.83>
- Tannenbaum, D., Ditto, P. H., & Pizarro, D. A. (2007). *Different moral values produce different judgments of intentional action* [Unpublished manuscript]. University of California-Irvine.
- Tetlock, P. E., Self, W. T., & Singh, R. (2010). The punitiveness paradox: When is external pressure exculpatory – And when a signal just to spread blame? *Journal of Experimental Social Psychology*, 46(2), 388395. <https://doi.org/10.1016/j.jesp.2009.11.013>
- Weiner, B. (1995). *Judgments of responsibility: A foundation for a theory of social conduct*. Guilford Press.
- Zucchelli, M. M., Starita, F., Bertini, C., Giusberti, F., & Ciaramelli, E. (2019). Intentionality attribution and emotion: The Knobe Effect in alexithymia. *Cognition*, 191, 103978. <https://doi.org/10.1016/j.cognition.2019.05.015>

Supplementary Materials

Supplementary Materials

Download: https://collabra.scholasticahq.com/article/128423-asymmetries-in-attributions-of-blame-and-praise-intent-and-causality-free-will-responsibility-and-the-side-effect-effect/attachment/261982.docx?auth_token=IXslCjM40zGct8eutozu

Peer Review Communication

Download: https://collabra.scholasticahq.com/article/128423-asymmetries-in-attributions-of-blame-and-praise-intent-and-causality-free-will-responsibility-and-the-side-effect-effect/attachment/261983.docx?auth_token=IXslCjM40zGct8eutozu

Asymmetries in attributions of blame and praise, intent, and causality: Free will, responsibility, and the side-effect effect

Supplementary

Power analyses.....	2
Study 1	11
Materials of Study 1	13
Manipulation Check.....	15
Dependent variables.....	15
Study 2	22
Materials of Study 2.....	22
Manipulation Check.....	24
Dependent variables.....	24
Additional results	26
References.....	39

Power analyses

We conducted a power analysis before collecting the data for the Study 2. The power analysis was based on the results of Study 1. Our aim with Study 2 was to detect the weakest effect reported in Study 1 at .95 power ($\alpha = 0.05$). The largest required sample based on the power analyses is **1086**.

Details:

Power analyses

The largest required sample based on the power analyses is **1086**.

Intent side-effect effect

In study 1, the side-effect effect for different universes produced following effect sizes:

Deterministic universe $d = 1.51$; indeterministic universe $d = 1.61$; control condition $d = 1.84$.

t tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size $|\rho| = 1.51$

α err prob = 0.05

Power ($1 - \beta$ err prob) = 0.95

Output: Noncentrality parameter $\delta = 3.5412639$

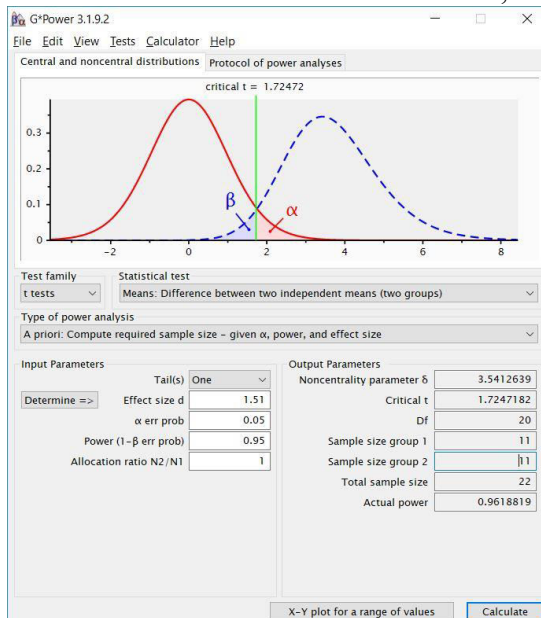
Critical $t = 1.7247182$

Df = 20

Total sample size = 22

Actual power = 0.9618819

Based on 1.51 as the lowest effect size, the required sample is **22**.



Causality side-effect effect

In study 1, the side-effect effect of causality different universes produced following effect sizes: Deterministic universe $d = 1.05$; indeterministic universe $d = 0.90$; control condition $d = 1.03$.

t tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size $|\rho| = 0.90$

α err prob = 0.05

Power ($1-\beta$ err prob) = 0.95

Output: Noncentrality parameter $\delta = 3.3674916$

Critical t = 1.6735649

Df = 54

Total sample size = 56

Actual power = 0.9535206

Based on 1.51 as the lowest effect size, the required sample is 56.

Free-will and blame association

In study 1, the correlation between free will attribution and blame we obtained the following Pearson correlation coefficient : $r = .532$. Using G*Power alpha = .05, two-tail (direction of hypothesis not determined) effect size $r = .532$ and power .95 we require a sample of 29.

t tests - Correlation: Point biserial model

Analysis: A priori: Compute required sample size

Input: Tail(s) = Two

Effect size $|\rho| = 0.532$

α err prob = 0.05

Power ($1-\beta$ err prob) = 0.95

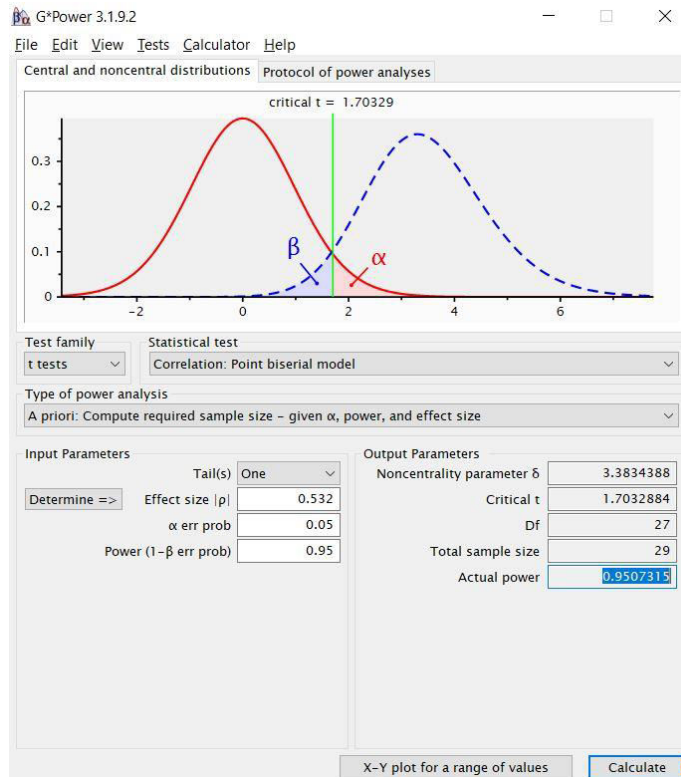
Output: Noncentrality parameter $\delta = 3.3834388$

Critical t = 1. 1.7032884

Df = 27

Total sample size = 29

Actual power = 0.9507315



In study 1, the partial correlation between free will attribution and blame after controlling for intention and casualty we obtained the following Pearson correlation coefficient: $r = .510$
 In study 1, the correlation between free will attribution and blame we obtained the following Pearson correlation coefficient: $r = .510$. Using G*Power alpha = .05, two-tail (direction of hypothesis not determined), $r = 0.510$ and power .95 we require a sample of **33**.

t tests - Correlation: Point biserial model

Analysis: A priori: Compute required sample size

Input: Tail(s) = Two

Effect size $|p| = 0.510$

α err prob = 0.05

Power (1- β err prob) = 0.95

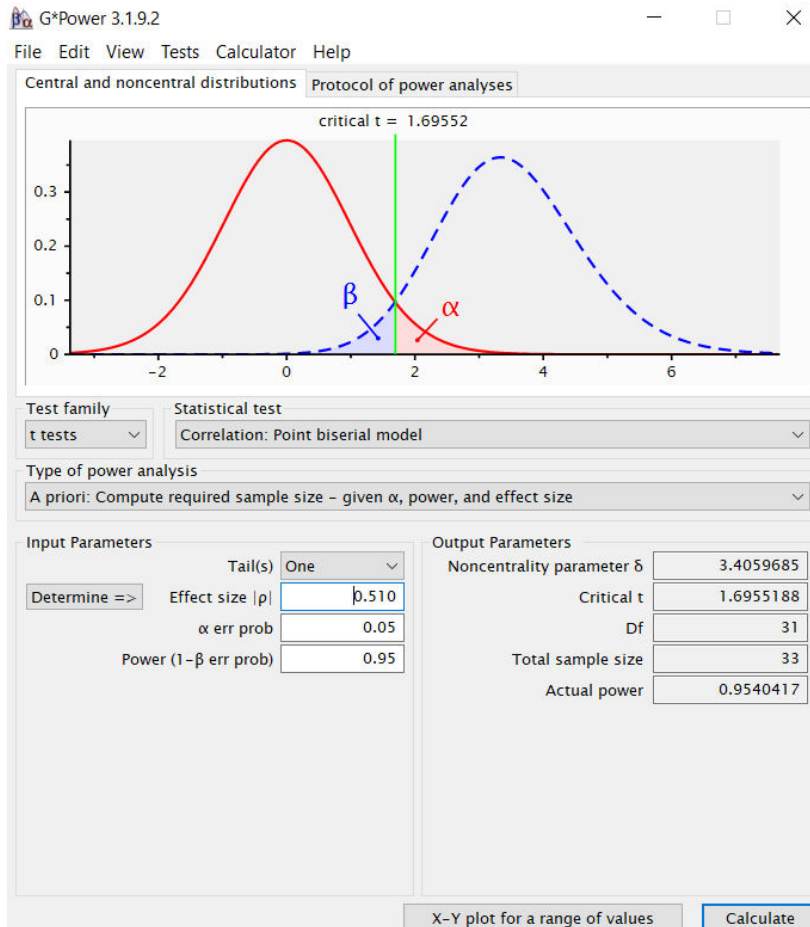
Output: Noncentrality parameter $\delta = 3.4059685$

Critical t = 1.6955188

Df = 31

Total sample size = 33

Actual power = 0.9540417



In study 1, the independent t-test between blame and praise produced following effect size: $d = 0.77$

t tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size $|\rho| = 0.7660409$

α err prob = 0.05

Power (1- β err prob) = 0.95

Output: Noncentrality parameter $\delta = 3.3563522$

Critical t = 1.6657069

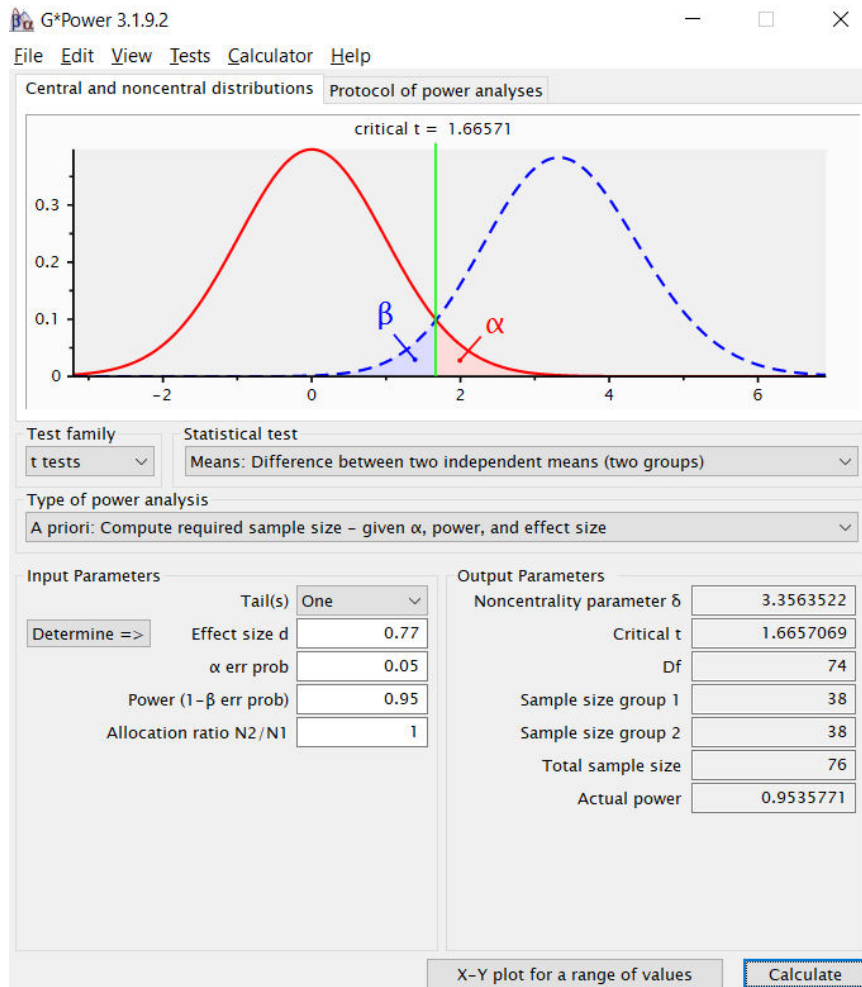
Df = 74

Sample size group 1 = 38

Sample size group 2 = 38

Actual power = 0.9535771

Based on 0.77 as the effect size, the required sample is 76.



In study 1, the independent t-test on attribution of blame for harmful outcome between indeterministic universe and deterministic universe produced following effect size: $d = -1.125$, with required sample of 36.

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size d = 1.125

α err prob = 0.05

Power (1- β err prob) = 0.95

Allocation ratio N2/N1 = 1

Output: Noncentrality parameter δ = 3.3750000

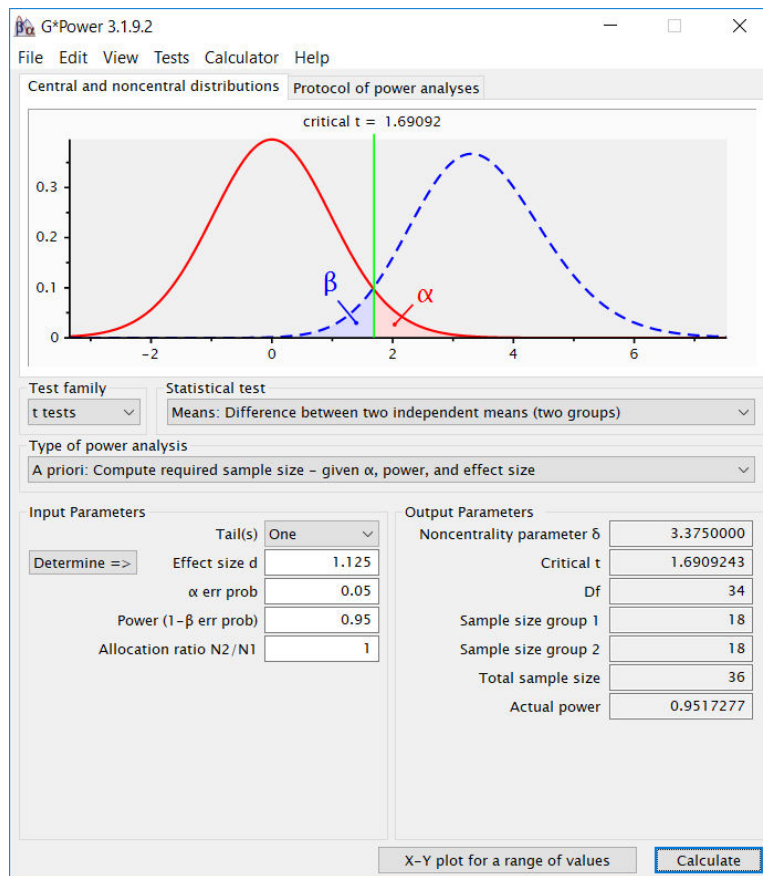
Critical t = 1.6909243

Df = 34

Sample size group 1 = 18

Sample size group 2 = 18

Total sample size = 36



Free-will and regret association

In a different project we measures free will and regret attributions and the correlational association between the two was $r = .14$, converted to Cohen's d is 0.28, which for power of 95% requires a sample size of **554**.

tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size d = 0.28

α err prob = 0.05

Power (1- β err prob) = 0.95

Allocation ratio N2/N1 = 1

Output: Noncentrality parameter δ = 3.2952086

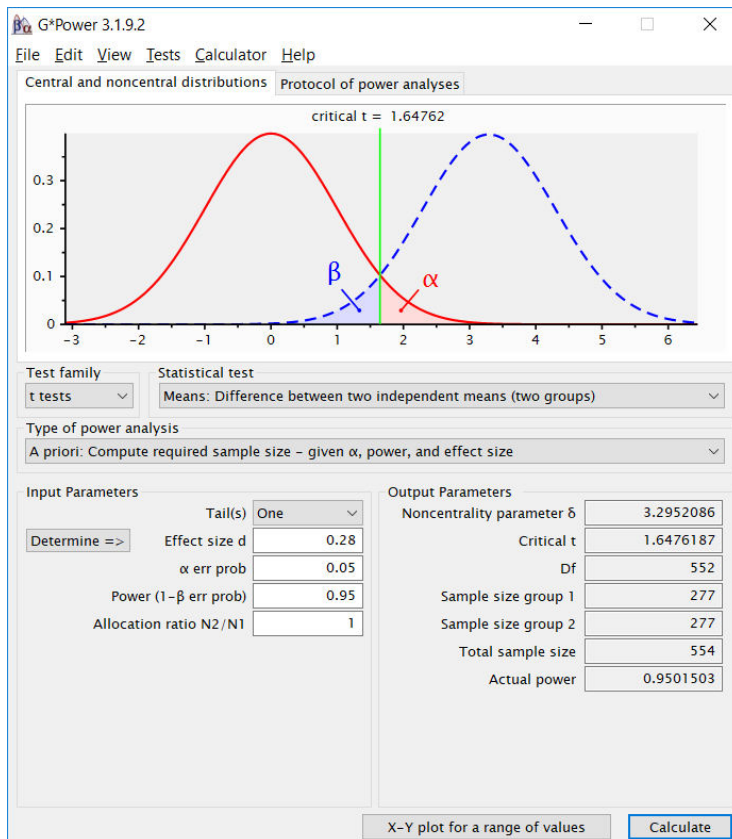
Critical t = 1.6476187

Df = 552

Sample size group 1 = 277

Sample size group 2 = 277

Total sample size = 554



Bad is stronger than good

In study 1, the independent t-test comparing attribution of blame for negative outcome condition and attribution of praise for positive outcomes produced following effect size: $d = 1.36$, with required sample of 26.

tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size d = 1.361847

α err prob = 0.05

Power (1- β err prob) = 0.95

Allocation ratio N2/N1 = 1

Output: Noncentrality parameter δ = 3.4720422

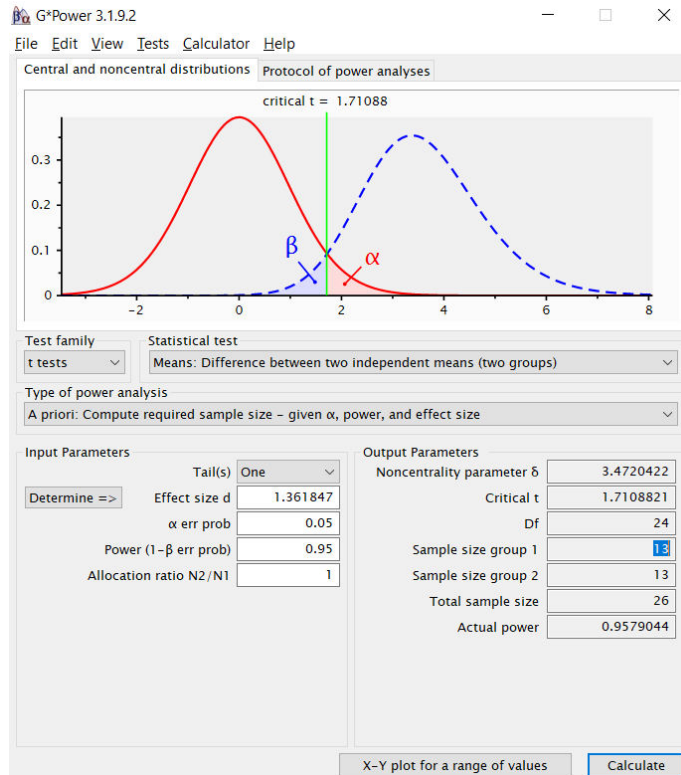
Critical t = 1.7108821

Df = 24

Sample size group 1 = 13

Sample size group 2 = 13

Total sample size = 26



Bad is free than good

In study 1, the independent t-test on attribution of free-will between harmful outcome between indeterministic universe and deterministic universe produced following effect size: $d = 0.199$, with required sample of **1086**.

tests - Means: Difference between two independent means (two groups)

Analysis: A priori: Compute required sample size

Input: Tail(s) = One

Effect size d = 0.1998839

α err prob = 0.05

Power (1-β err prob) = 0.95

Allocation ratio $N2/N1$ = 1

Output: Noncentrality parameter δ = 3.2935384

Critical t = 1.6462605

Df = 1084

Sample size group 1 = 543

Sample size group 2 = 543

Total sample size = 1086

Free will attributions with nuanced comparison

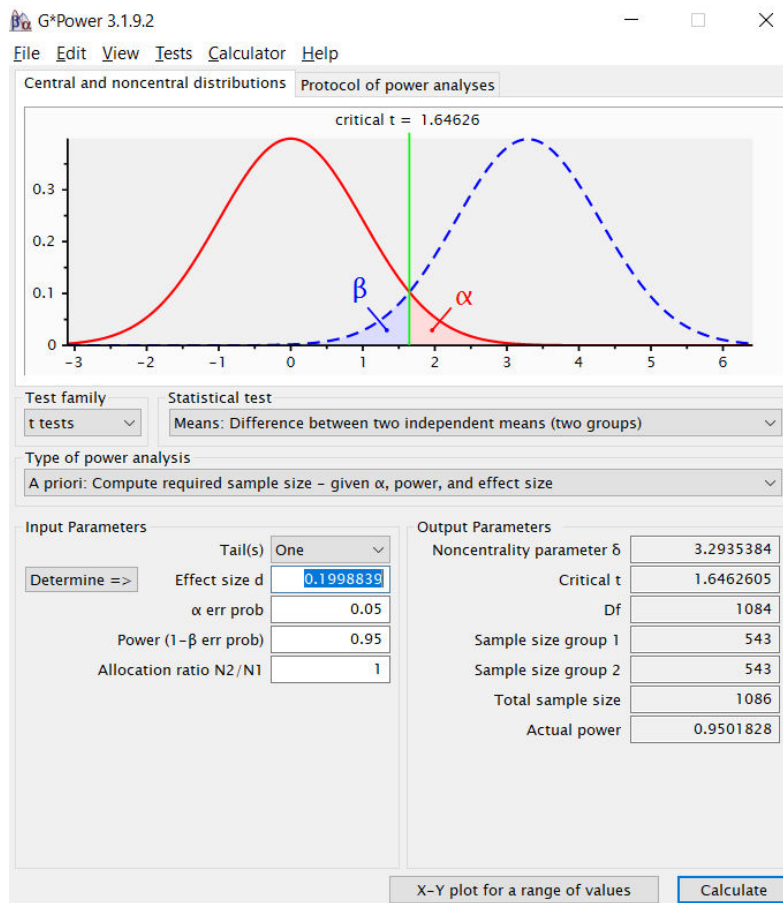
All universes Harm vs. Help; Cohen's $d = 0.199$; Required $N = 1086$ (shown above)

Deterministic universe - harm & help; Cohen's $d = 0.201337$; Required $N = 1070$

Indeterministic universe - harm & help; Cohen's $d = 0.259705$; Required $N = 644$

Unknown universe - harm & help; Cohen's $d = 0.492977$; Required $N = 180$

Control universe - harm & help; Cohen's $d = 0.236293$; Required $N = 778$



Study 1

Important note

The analysis presented in Study 1 is based on the data collected in testing another set of published hypotheses by Feldman and Chandrashekar (2018). In Feldman and Chandrashekar's (2018) study, the key measures of interest were not the same as the one presented here, however, because the experimental manipulations were identical, we included the measures of interest in the Feldman and Chandrashekar (2018). The commonality between Feldman and Chandrashekar (2018) and Study 1 is the manipulation of the universe conditions. The study design of Feldman and Chandrashekar's (2018) included four between-subjects universe manipulations (i.e., "deterministic universe," "indeterministic universe," "uncertain universe," "Control condition").

For the present investigation, we only exclude the responses from the participants assigned to the uncertain universe (in which it is unclear to agents whether human behavior is determined or undetermined), as this manipulation is not relevant for the current set of theoretical predictions. As part of the experimental materials of Study 1 we document all the procedures, including that of the universe manipulations noted in the (Feldman and Chandrashekar (2018).

Table S1
Experimental Design of Study 1

Study 1				
Participants were randomly assigned to 1 out of 8 different default and framing conditions and were required to confirm their choices accordingly. Experimental conditions varied in the structure of the question, i.e., the structure of the question (the DV) presented to the participants at the end of the health survey varied on framing and defaults.				
Independent Variable 1: Universe conditions	IV condition 1: Deterministic	IV condition 2: Indeterministic	IV condition 3: Control condition	IV condition 4: Uncertain universe
Independent Variable 1: Universe conditions	Outcome: Help		Outcome: Harm	
Dependent Variables	<ul style="list-style-type: none">• Attributions of blame• Attributions of intentionality• Attributions of causality• Attributions of freewill• Attributions of praise			

Note. Responses recorded as part of "Uncertain universe" condition (which was integral part of Feldman and Chandrashekar, 2018) were not part of Study 1.

Materials of Study 1

'Determinism universe - harm the environment' condition

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Determinism universe - help the environment' condition

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

'Indeterministic universe - harm the environment' condition

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just

want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Indeterministic universe - help the environment' condition

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

'Unknown universe - harm the environment' condition

Imagine a universe (Universe D) in which it is possible that everything that happens is completely caused by whatever happened before it. But in this universe, it is unclear whether human action follows this rule or if it is an exception to this rule.

For example, one-day John decided to have French Fries at lunch. In this universe, it is unclear whether John's decision in this universe was or was not completely caused by what happened before it. Assuming everything in the universe was exactly the same up until John made his decision, it is unclear whether or not John could have decided to not have French Fries, and whether he could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Unknown universe - help the environment' condition

Imagine a universe (Universe D) in which it is possible that everything that happens is completely caused by whatever happened before it. But in this universe, it is unclear whether human action follows this rule or if it is an exception to this rule.

For example, one-day John decided to have French Fries at lunch. In this universe, it is unclear whether John's decision in this universe was or was not completely caused by what happened before it. Assuming everything in the universe was exactly the same up until John made his decision, it is unclear whether or not John could have decided to not have French Fries, and whether he could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

'Control universe - harm the environment' condition

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Control universe - help the environment'. condition

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

Manipulation Check

To make sure you understood the scenario - what was the environmental outcome of the chairman's decision to start the new program?:

- (a) The environment was helped.
- (b) The environment was harmed.
- (c) The scenario doesn't say.

Dependent variables

Attributions of intentionality

1. Do you agree with the following statement? - In Universe D, the chairman intentions were to have such implications of the new program on the environment?
2. Do you agree with the following statement? - In Universe D, did the chairman intentionally affect the environment?

Attributions of causality

1. Do you agree with the following statement? - In Universe D, the chairman knew the implications of the new program on the environment? (1 = strongly disagree; 6 = strongly agree)

2. Do you agree with the following statement? - In Universe D, the chairman understood the implications of the new program on the environment? (1 = strongly disagree; 6 = strongly agree)

Attributions of freewill

1. Do you agree with the following statement? - In Universe D, the chairman had to choose what he chose, and could not have chosen to do otherwise. (R) (1 = strongly disagree; 6 = strongly agree)
2. Do you agree with the following statement? - In Universe D, the chairman was free to choose not to start the new program. (1 = strongly disagree; 6 = strongly agree)

Attributions of blame

1. Do you agree with the following statement? - In Universe D, the chairman should be criticized for his actions if they led to negative outcomes. (1 = strongly disagree; 6 = strongly agree)

Attributions of praise

1. Do you agree with the following statement? - In Universe D, the chairman should be applauded for his actions if they led to positive outcomes. (1 = strongly disagree; 6 = strongly agree)

Table S2

Study 1: Results of 2x2 ANOVA testing the effects of type of outcome and type of universe on attributions of free will, Intentionality, and Causality.

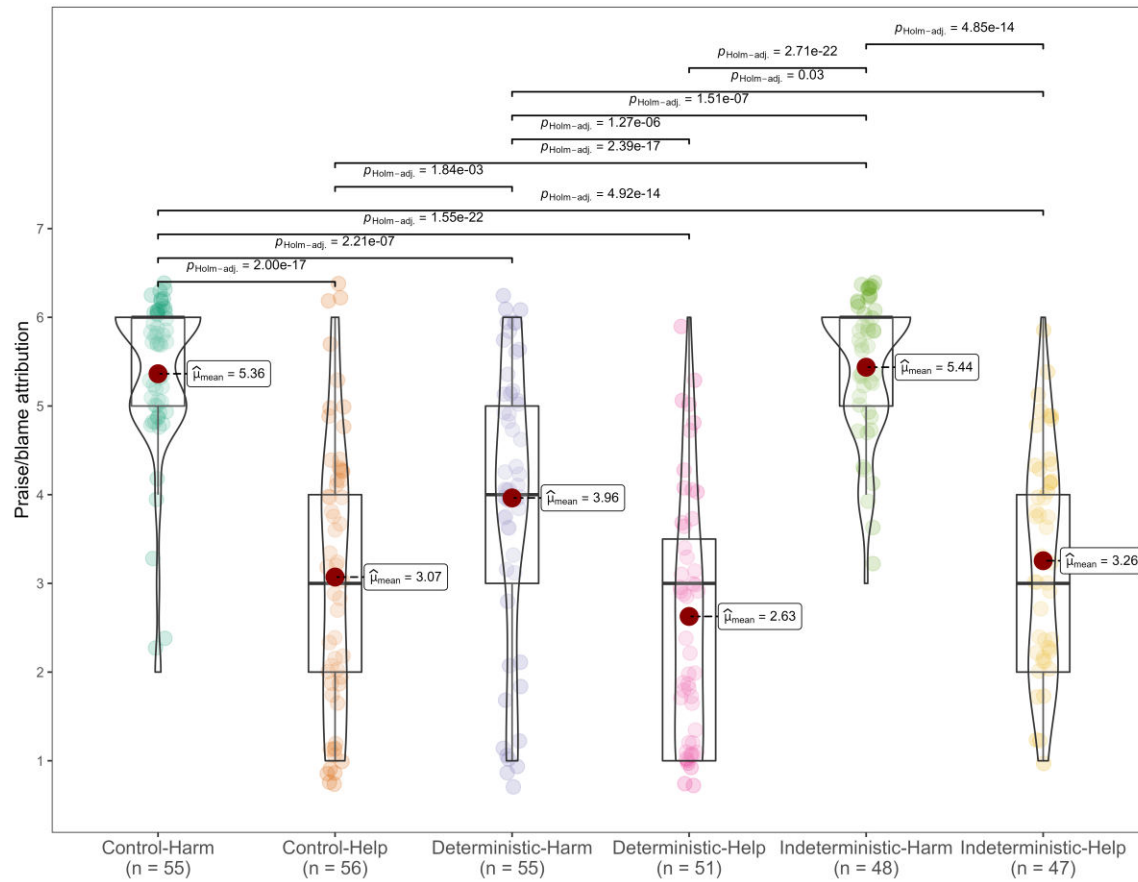
Factor	Praise/Blame attribution					Intentionality attribution					Causality attribution					Free Will				
	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p
Type of outcome (Help vs Harm)	91.77	1	154.93	<.001	0.32	97.25	1	157.07	<.001	0.33	29.95	1	40.28	<.001	0.13	2.21	1	3.96	.14	.01
Type of universe	32.75	1	55.29	<.001	0.14	4.21	1	6.79	.042	.02	4.11	1	5.52	.044	.02	183.34	1	328.58	<.001	.48
Type of outcome × Type of universe	5.31	1	8.96	.022	0.03	1.61	1	2.60	.206	0.01	0.09	1	0.12	.768	0.00	0.12	1	0.21	.73	.001

Note. *df* = degree of freedom, *MS* = Mean Sum of Squares.

Figure S1

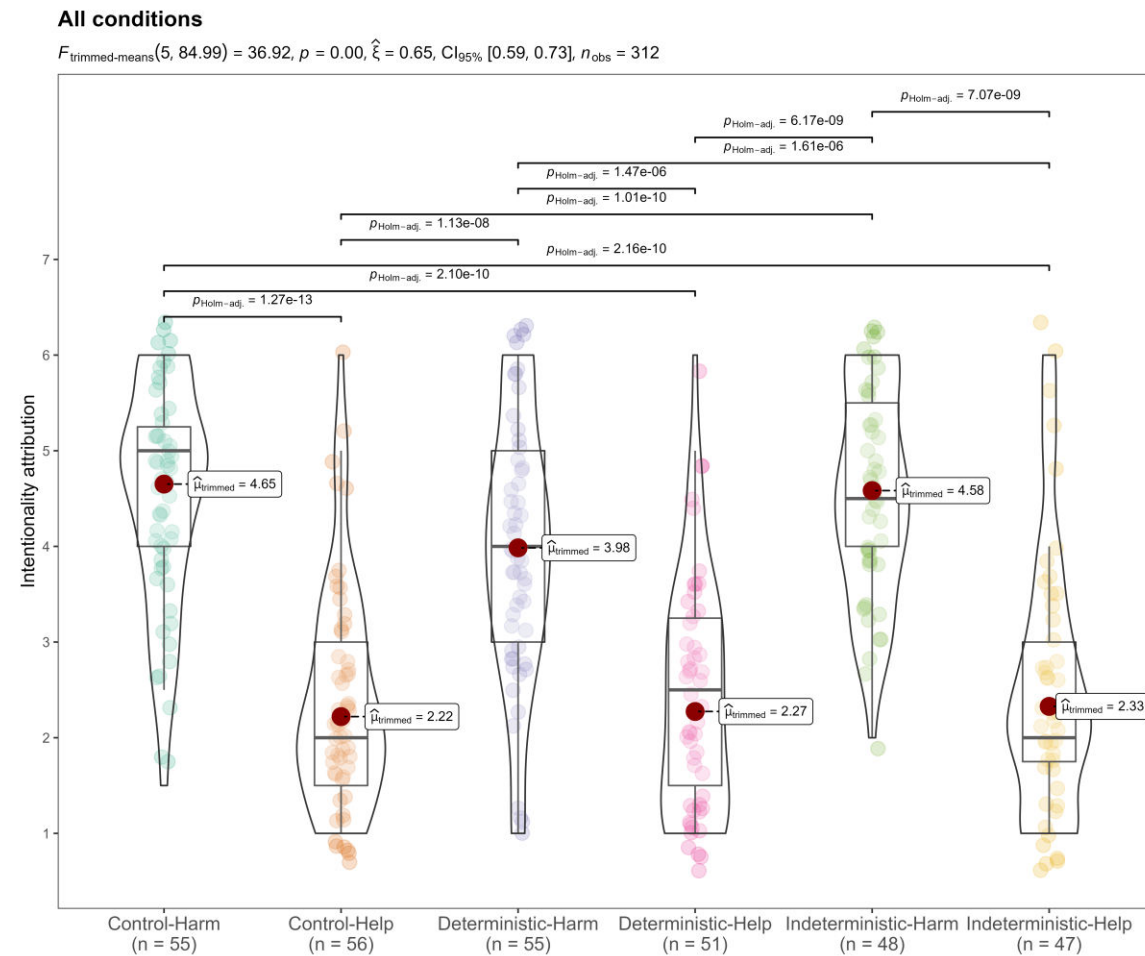
*Study 1 Attribution of Praise/blame across all conditions***All conditions**

$F_{\text{Fisher}}(5, 306) = 44.81, p = 1.26\text{e-}34, \omega_p^2 = 0.41, \text{CI}_{95\%} [0.34, 1.00], n_{\text{obs}} = 312$



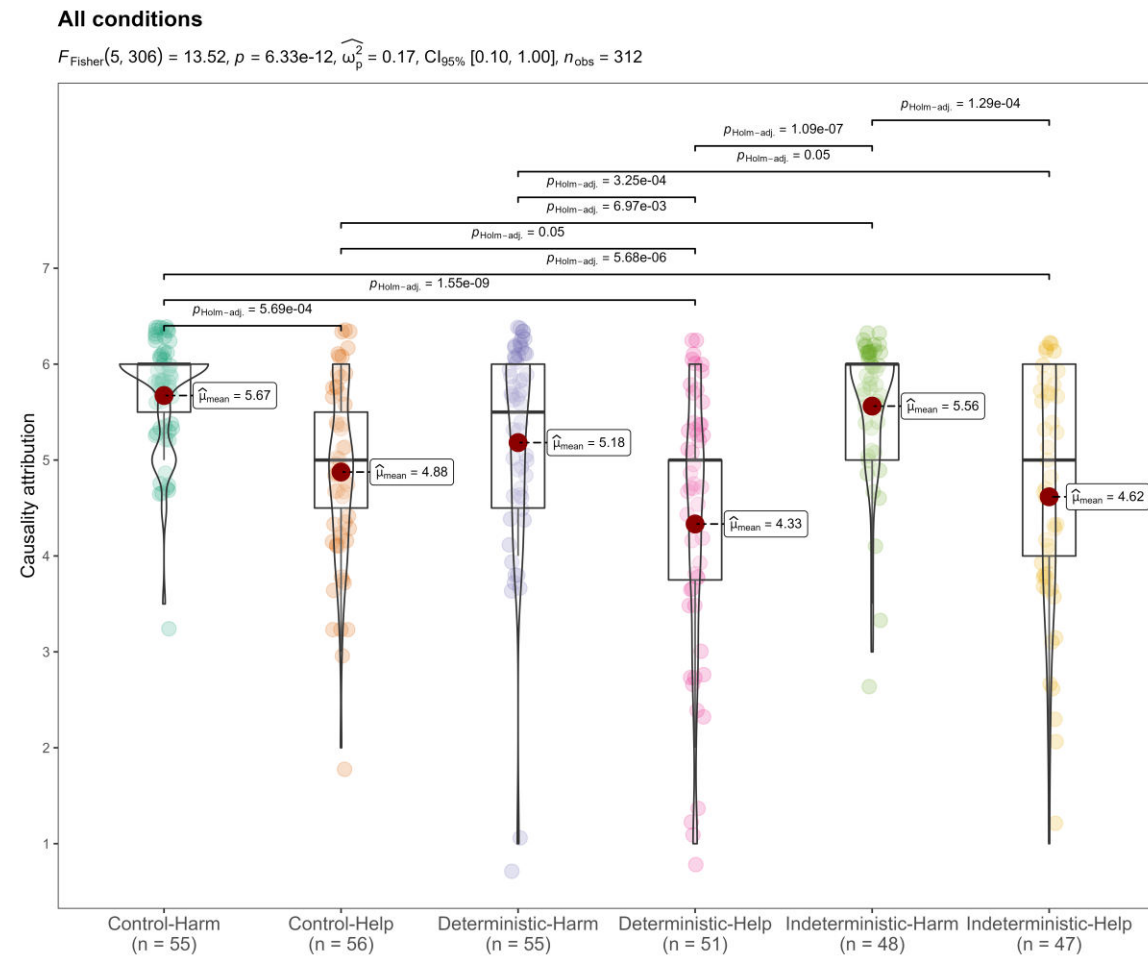
Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S2

Study 1 Attribution of intentionality across all conditions

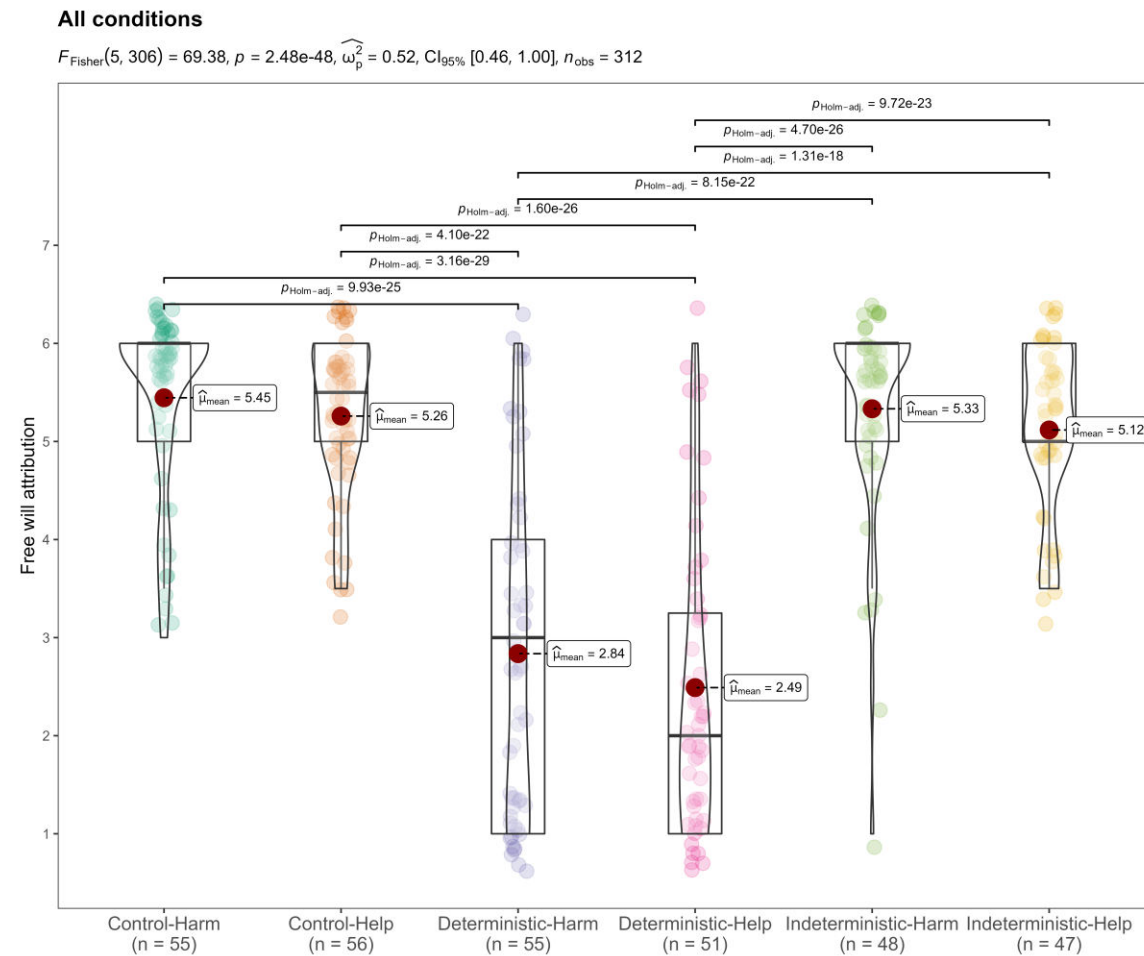
Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S3

Study 1 Attribution of causality across all conditions

Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S4

Study 1 Attribution of free will across all conditions

Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Study 2

Table S3

Experimental Design of Study 2

<p>Study 2</p> <p>Participants were randomly assigned to 1 out of 6 different default and framing conditions and were required to confirm their choices accordingly. Experimental conditions varied in the structure of the question, i.e., the structure of the question (the DV) presented to the participants at the end of the health survey varied on framing and defaults.</p>			
Independent Variable 1: Universe conditions	IV condition 1: Deterministic	IV condition 2: Indeterministic	IV condition 3: Control condition
Independent Variable 1: Universe conditions	Outcome: Help		Outcome: Harm
Dependent Variables	<ul style="list-style-type: none"> • Attributions of blame • Attributions of intentionality • Attributions of causality • Attributions of freewill • Attributions of praise • Attributions of regret 		

Materials of Study 2

'Determinism universe - harm the environment' condition

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Determinism universe - help the environment' condition

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

'Indeterministic universe - harm the environment' condition

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment". The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Indeterministic universe - help the environment' condition

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

In Universe D there is a company.

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment". The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

'Control universe - harm the environment' condition

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, but it will also harm the environment".

The chairman of the board answered, "I don't care at all about harming the environment. I just want to make as much profit as I can. Let's start the new program." They started the new program. Sure enough, the environment was harmed.

'Control universe - help the environment'. condition

The vice-president of this company went to the chairman of the board and said, "We are thinking of starting a new program. It will help us increase profits, and it will also help the environment".

The chairman of the board answered, "I don't care at all about helping the environment. I just want to make as much profit as I can. Let's start the new program". They started the new program. Sure enough, the environment was helped.

Manipulation Check

To make sure you understood the scenario - what was the environmental outcome of the chairman's decision to start the new program?:

- (a) The environment was helped.
- (b) The environment was harmed.
- (c) The scenario doesn't say.

Dependent variables

Attributions of intentionality

- 1. Do you agree with the following statement? - In Universe D, the chairman intentions were to have such implications of the new program on the environment?
- 2. Do you agree with the following statement? - In Universe D, did the chairman intentionally affect the environment?

Attributions of causality

- 1. Do you agree with the following statement? - In Universe D, the chairman knew the implications of the new program on the environment? (1 = strongly disagree; 6 = strongly agree)
- 2. Do you agree with the following statement? - In Universe D, the chairman understood the implications of the new program on the environment? (1 = strongly disagree; 6 = strongly agree)

Attributions of freewill

- 1. Do you agree with the following statement? - In Universe D, the chairman had to choose what he chose, and could not have chosen to do otherwise. (R) (1 = strongly disagree; 6 = strongly agree)
- 2. Do you agree with the following statement? - In Universe D, the chairman was free to choose not to start the new program. (1 = strongly disagree; 6 = strongly agree)

Attributions of blame

- 1. Do you agree with the following statement? - In Universe D, the chairman should be criticized for his actions if they led to negative outcomes. (1 = strongly disagree; 6 = strongly agree)

Attributions of praise

1. Do you agree with the following statement? - In Universe D, the chairman should be applauded for his actions if they led to positive outcomes. (1 = strongly disagree; 6 = strongly agree)

Attributions of regret

1. Do you agree with the following statement? - In Universe D, the chairman would regret his decision if he learned that his actions led to the environment being harmed. (1 = strongly disagree; 6 = strongly agree)

Additional results

Table S4

Descriptive statistics grouped by experimental conditions

Study	Condition	Dimension	n	Mean	SD	Median	Skew	Kurtosis
Study 1	Control	Praise	111	2.78	1.58	2.00	0.44	-1.08
		Blame	111	4.99	1.14	5.00	-1.3	1.3
		Intention	111	3.45	1.56	3.50	0.08	-1.29
		Freewill	111	5.35	0.84	5.50	-1.18	0.24
		Causality	111	5.27	0.84	5.50	-1.1	1.04
	Deterministic Universe	Praise	106	2.75	1.39	3.00	0.27	-1.07
		Blame	106	3.72	1.48	4.00	-0.31	-0.75
		Intention	106	3.2	1.52	3.00	0.18	-0.9
		Freewill	106	2.67	1.61	2.00	0.64	-0.82
		Causality	106	4.77	1.34	5.00	-1.29	1.26
	Indeterministic Universe	Praise	95	3.4	1.51	4.00	-0.07	-1.14
		Blame	95	4.91	1.24	5.00	-1.39	1.65
		Intention	95	3.55	1.58	3.50	0.02	-1.14
		Freewill	95	5.23	0.96	5.50	-1.58	3.04
		Causality	95	5.09	1.11	5.50	-1.23	1.07
Study 2	Control	Praise	367	2.64	1.52	2.00	0.58	-0.76
		Blame	367	5.01	1.19	5.00	-1.31	1.20
		Intention	367	3.44	1.64	3.50	0.05	-1.23
		Freewill	367	5.16	0.98	5.50	-1.23	1.31
		Causality	367	5.08	1.13	5.00	-1.50	2.10
		Regret	367	2.16	1.34	2.00	1.24	0.80
	Deterministic Universe	Praise	359	3.01	1.54	3.00	0.31	-0.99
		Blame	359	3.82	1.65	4.00	-0.29	-1.16
		Intention	359	3.31	1.60	3.50	0.13	-1.15
		Freewill	359	2.10	1.27	2.00	1.21	0.76
		Causality	359	5.09	1.06	5.00	-1.38	1.87
		Regret	359	2.29	1.36	2.00	1.08	0.32
	Indeterministic Universe	Praise	367	3.20	1.64	3.00	0.16	-1.19
		Blame	367	4.98	1.27	5.00	-1.41	1.42
		Intention	367	3.44	1.60	3.50	-0.02	-1.21
		Freewill	367	5.31	0.92	6.00	-1.43	1.61
		Causality	367	5.07	1.10	5.00	-1.45	1.96
		Regret	367	2.20	1.27	2.00	1.18	0.85

Table S5

Study 1 Descriptive statistics grouped by experimental conditions And outcome of the scenario

Experimental condition	Outcome	Dimension	<i>n</i>	Mean	SD	Median	Skew	Kurtosis
Control	Harm	Praise	55	2.49	1.60	2.00	0.74	-0.90
		Blame	55	5.36	0.93	6.00	-1.98	4.27
		Intention	55	4.52	1.14	5.00	-0.63	-0.36
		Freewill	55	5.45	0.91	6.00	-1.44	0.68
		Causality	55	5.67	0.55	6.00	-1.70	2.83
	Help	Praise	56	3.07	1.52	3.00	0.18	-1.03
		Blame	56	4.63	1.21	5.00	-0.89	0.30
		Intention	56	2.40	1.16	2.00	1.03	0.70
		Freewill	56	5.26	0.77	5.50	-0.88	-0.25
		Causality	56	4.88	0.89	5.00	-0.63	0.42
Deterministic Universe	Harm	Praise	55	2.85	1.41	3.00	0.10	-1.30
		Blame	55	3.96	1.57	4.00	-0.56	-0.67
		Intention	55	3.95	1.38	4.00	-0.32	-0.43
		Freewill	55	2.84	1.71	3.00	0.43	-1.13
		Causality	55	5.18	1.12	5.50	-1.92	4.52
	Help	Praise	51	2.63	1.39	3.00	0.45	-0.80
		Blame	51	3.45	1.35	4.00	-0.11	-0.70
		Intention	51	2.40	1.22	2.50	0.69	-0.08
		Freewill	51	2.49	1.49	2.00	0.87	-0.43
		Causality	51	4.33	1.43	5.00	-0.87	-0.02
Indeterministic Universe	Harm	Praise	48	3.54	1.71	4.00	-0.21	-1.34
		Blame	48	5.44	0.77	6.00	-1.17	0.57
		Intention	48	4.54	1.09	4.50	-0.17	-0.94
		Freewill	48	5.33	1.06	6.00	-2.07	4.63
		Causality	48	5.56	0.70	6.00	-1.84	3.23
	Help	Praise	47	3.26	1.28	3.00	0.02	-1.09
		Blame	47	4.36	1.39	5.00	-0.98	0.14
		Intention	47	2.54	1.35	2.00	0.94	0.23
		Freewill	47	5.12	0.84	5.00	-0.72	-0.62
		Causality	47	4.62	1.25	5.00	-0.63	-0.14

Table S6

Study 2 Descriptive statistics grouped by experimental conditions And outcome of the scenario

Experimental condition	Outcome	Dimension	<i>n</i>	Mean	SD	Median	Skew	Kurtosis
Control	Help	Praise	185	2.16	1.54	1.00	1.22	0.29
		Blame	185	5.47	0.92	6.00	-2.30	6.24
		Intention	185	4.49	1.30	4.50	-0.62	-0.31
		Freewill	185	5.48	0.88	6.00	-1.98	3.99
		Causality	185	5.48	0.88	6.00	-2.46	7.56
		Regret	185	2.09	1.39	2.00	1.33	0.95
	Harm	Praise	182	3.13	1.33	3.00	0.17	-0.70
		Blame	182	4.54	1.26	5.00	-0.82	0.00
		Intention	182	2.38	1.21	2.00	0.63	-0.69
		Freewill	182	4.85	0.97	5.00	-0.86	0.74
		Causality	182	4.68	1.22	5.00	-1.02	0.62
		Regret	182	2.23	1.29	2.00	1.14	0.62
Deterministic Universe	Help	Praise	181	3.08	1.63	3.00	0.24	-1.17
		Blame	181	4.15	1.63	5.00	-0.52	-0.99
		Intention	181	4.22	1.41	4.50	-0.49	-0.60
		Freewill	181	2.17	1.35	2.00	1.22	0.59
		Causality	181	5.42	0.82	6.00	-2.06	5.84
		Regret	181	2.40	1.44	2.00	0.95	-0.01
	Harm	Praise	178	2.94	1.45	3.00	0.35	-0.81
		Blame	178	3.49	1.60	4.00	-0.11	-1.20
		Intention	178	2.40	1.23	2.00	0.69	-0.32
		Freewill	178	2.03	1.18	2.00	1.12	0.64
		Causality	178	4.76	1.16	5.00	-0.93	0.51
		Regret	178	2.17	1.27	2.00	1.19	0.65
Indeterministic Universe	Help	Praise	183	3.50	1.83	4.00	-0.12	-1.46
		Blame	183	5.42	0.79	6.00	-1.55	2.86
		Intention	183	4.52	1.15	5.00	-0.63	-0.17
		Freewill	183	5.42	0.82	6.00	-1.52	1.96
		Causality	183	5.48	0.69	6.00	-1.56	3.20
		Regret	183	2.20	1.32	2.00	1.25	0.88
	Harm	Praise	184	2.91	1.37	3.00	0.35	-0.67
		Blame	184	4.53	1.49	5.00	-0.89	-0.25
		Intention	184	2.36	1.22	2.00	0.73	-0.14
		Freewill	184	5.20	1.01	5.50	-1.28	1.03
		Causality	184	4.67	1.28	5.00	-0.95	0.30
		Regret	184	2.20	1.23	2.00	1.08	0.72

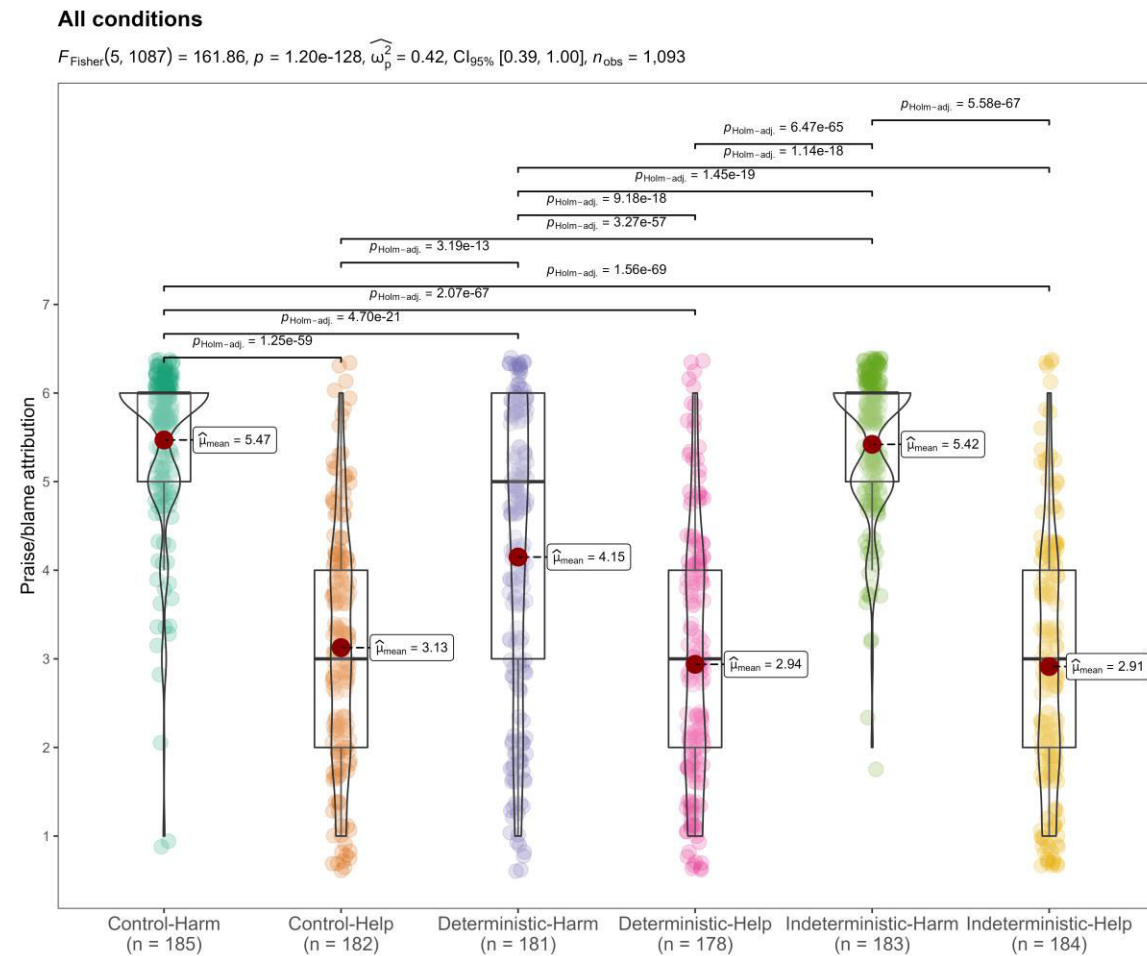
Table S7

Study 2 full results of 2x2 ANOVA testing the effects of type of outcome and type of universe on attributions of free will, intentionality, and causality.

Factor	Praise/Blame attribution					Intentionality attribution					Causality attribution					Free Will attribution				
	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p	<i>F</i>	<i>df</i>	<i>MS</i>	<i>p</i>	η^2p
Type of outcome (Help vs Harm)	346.2	1	627.37	<.001	0.32	453.64	1	715.60	<.001	0.39	95.44	1	98.56	<.001	0.12	5.06	1	6.16	.03	0.00
Type of universe	38.9	1	70.48	<.001	.05	2.00	1	3.15	.16	.003	0.04	1	0.04	.85	0.00	1537.25	1	1871.40	<.001	0.00
Type of outcome x Type of universe	42.1	1	76.29	<.001	0.06	3.29	1	5.19	.070	0.005	0.93	1	0.96	.34	0.001	0.22	1	0.27	.64	0.00

Note. *df* = degree of freedom, *MS* = Mean Sum of Squares.

Figure S5

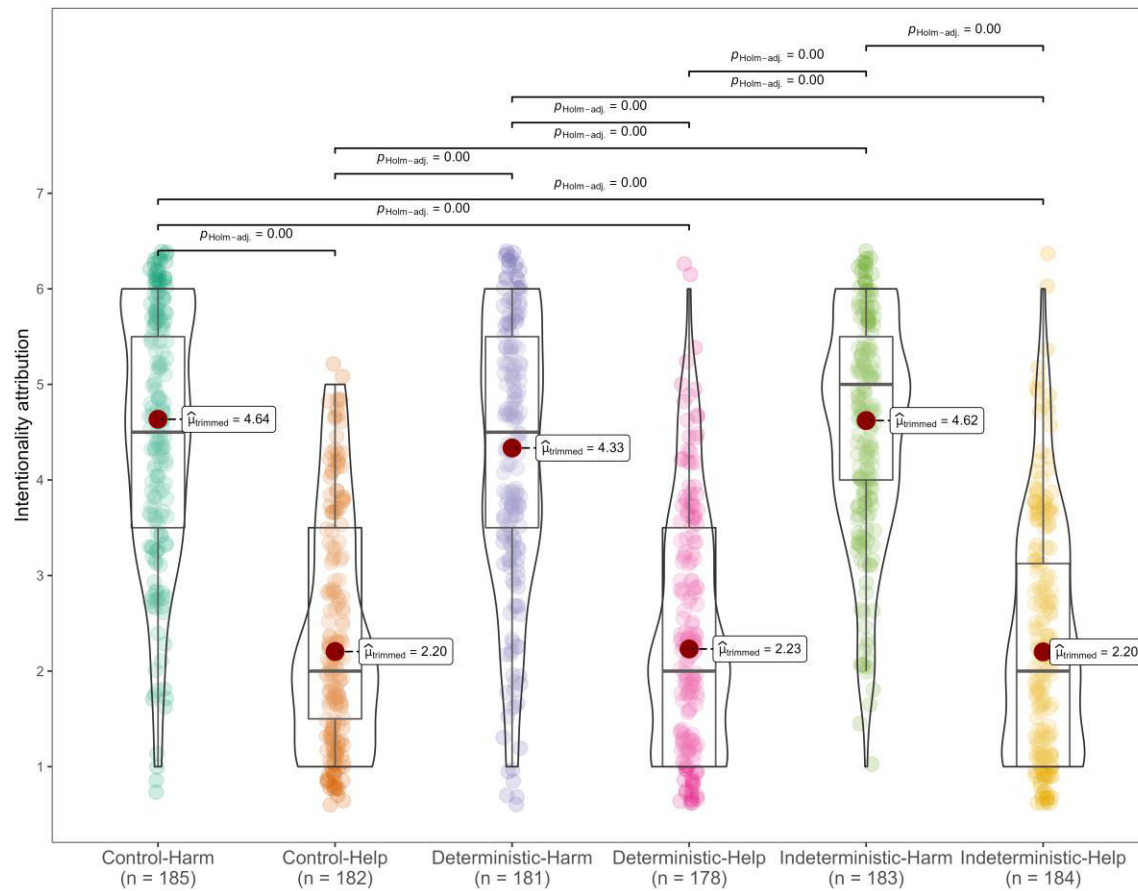
Study 2 Attribution of Praise/Blame across all conditions

Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S6

*Study 2 Attribution of intentionality across all conditions***All conditions**

$F_{\text{trimmed-means}}(5, 305.05) = 126.73, p = 0.00, \hat{\xi} = 0.67, CI_{95\%} [0.63, 0.71], n_{\text{obs}} = 1,093$

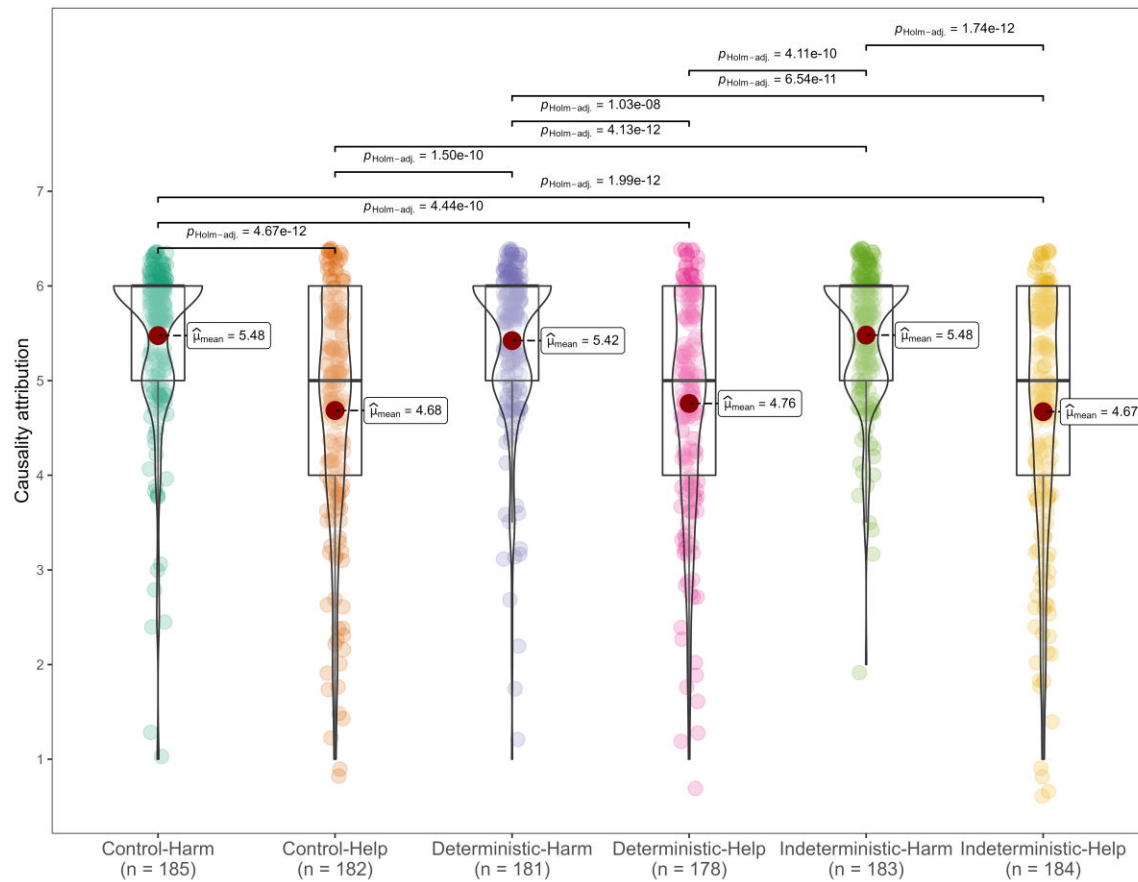


Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S7

*Study 2 Attribution of causality across all conditions***All conditions**

$F_{\text{Fisher}}(5, 1087) = 29.56, p = 3.26\text{e-}28, \hat{\omega}_p^2 = 0.12, \text{CI}_{95\%} [0.08, 1.00], n_{\text{obs}} = 1,093$

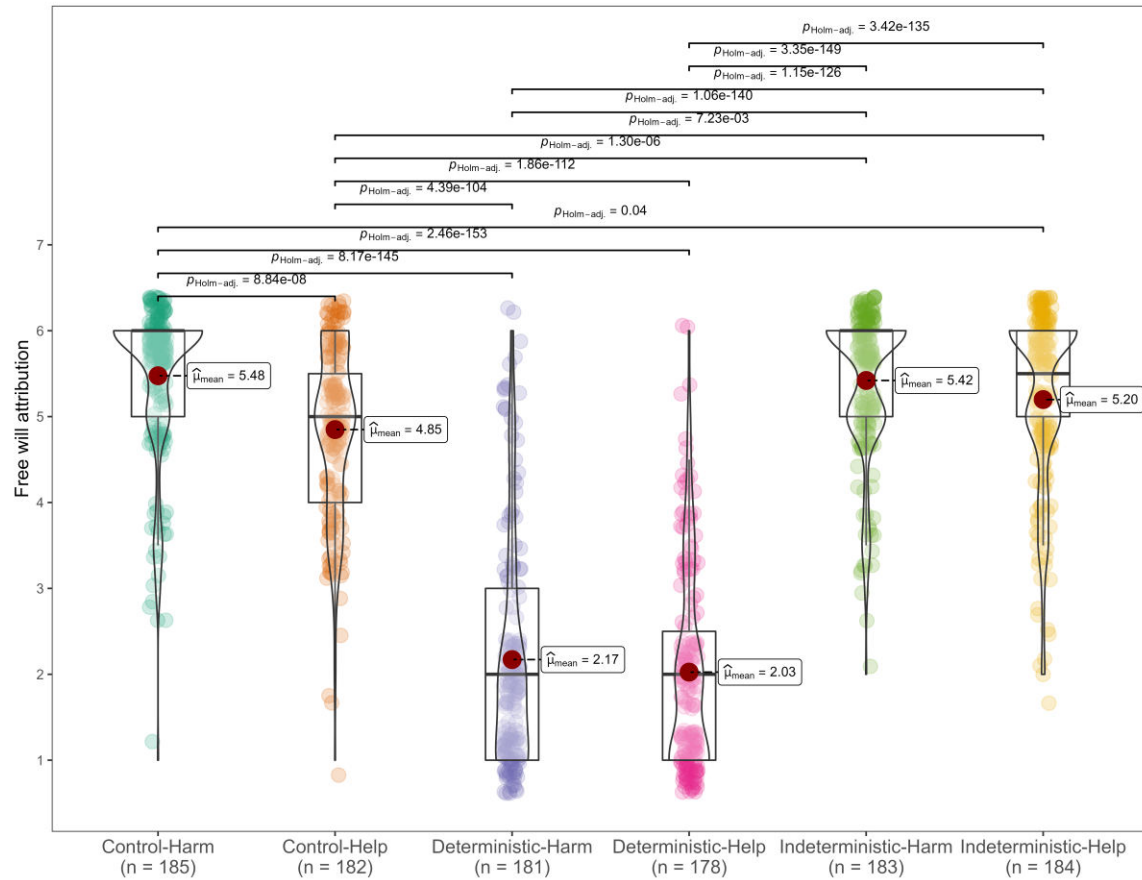


Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S8

*Study 2 Attribution of free will across all conditions***All conditions**

$F_{\text{Fisher}}(5, 1087) = 441.77, p = 7.85\text{e-}259, \hat{\omega}_p^2 = 0.67, \text{CI}_{95\%} [0.64, 1.00], n_{\text{obs}} = 1,093$

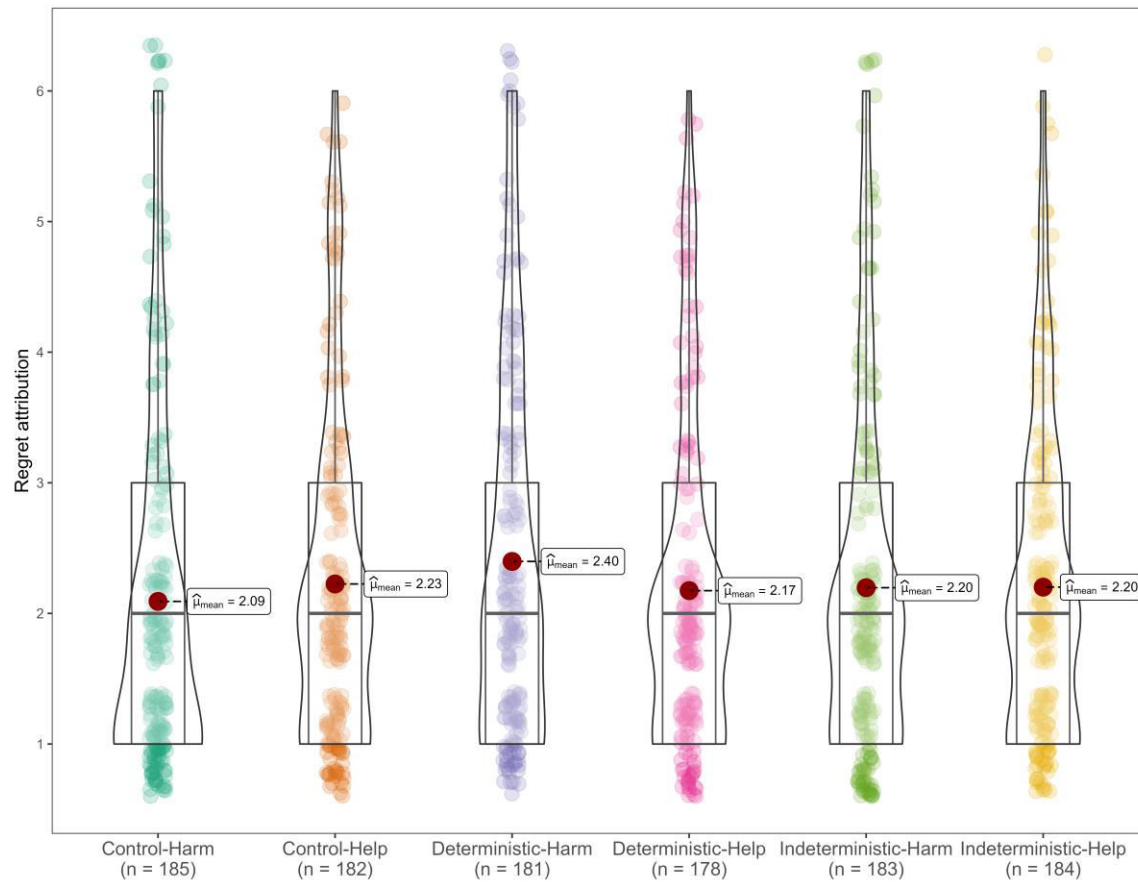


Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Figure S9

*Study 2 Attribution of regret across all conditions***All conditions**

$F_{\text{Fisher}}(5, 1087) = 1.05, p = 0.38, \hat{\omega}_p^2 = 2.47\text{e-}04, \text{CI}_{95\%} [0.00, 1.00], n_{\text{obs}} = 1,093$



Note. Violin plots displaying the distribution of responses, boxplots displaying the median, first, and third quartiles, and the red circle identifying the mean value.

Table S8
Study 1 correlations across all conditions

Type of universe	Variable	<i>M</i>	<i>SD</i>	1	2	3	4
Control (No universe, <i>n</i> =) 111)	1. Free will attributions	5.35	0.84				
	2. Intent attributions	3.43	1.55	.07 [-.12, .25]			
	3. Causality attributions	5.29	0.78	.24* [.06, .41]	.43** [.26, .57]		
	4. Praise attributions	2.82	1.58	-.08 [-.26, .11]	-.03 [-.22, .16]	-.07 [-.25, .12]	
	5. Blame attributions	5.03	1.10	.36** [.19, .50]	.36** [.19, .51]	.36** [.18, .51]	-.07 [-.25, .12]
Deterministic universe (<i>n</i> = 106)	1. Free will attributions	2.65	1.61				
	2. Intent attributions	3.20	1.52	.22* [.03, .39]			
	3. Causality attributions	4.78	1.35	-.19* [-.37, -.00]	.32** [.14, .48]		
	4. Praise attributions	2.75	1.40	.19* [.00, .37]	.31** [.13, .47]	-.02 [-.21, .17]	
	5. Blame attributions	3.73	1.48	.47** [.31, .61]	.27** [.09, .44]	.04 [-.15, .23]	.27** [.08, .44]
Indeterministic universe (<i>n</i> = 95)	1. Free will attributions	5.26	0.94				
	2. Intent attributions	3.55	1.60	.06 [-.14, .26]			
	3. Causality attributions	5.10	1.11	.30** [.10, .47]	.31** [.12, .48]		
	4. Praise attributions	3.39	1.52	.00 [-.20, .20]	.14 [-.07, .33]	-.00 [-.20, .20]	
	5. Blame attributions	4.91	1.25	.23* [.02, .41]	.32** [.13, .49]	.55** [.39, .67]	-.11 [-.31, .09]

Table S9
Study 2 correlations across all conditions

Type of universe	Variable	<i>M</i>	<i>SD</i>	1	2	3	4
Control (No universe, <i>n</i> = 358)	1. Free will attributions	5.16	0.98				
	2. Intent attributions	3.44	1.64	.12* [.02, .22]			
	3. Causality attributions	5.08	1.13	.34** [.25, .43]	.33** [.24, .42]		
	4. Praise attributions	2.64	1.52	-.28** [-.38, -.18]	-.19** [-.29, -.09]	-.24** [-.33, -.14]	
	5. Blame attributions	5.01	1.19	.47** [.39, .55]	.26** [.16, .35]	.27** [.17, .36]	-.32** [-.41, -.22]
Deterministic universe (<i>n</i> = 355)	1. Free will attributions	2.09	1.27				
	2. Intent attributions	3.31	1.61	.17** [.07, .27]			
	3. Causality attributions	5.11	1.05	-.13* [-.23, -.02]	.24** [.14, .34]		
	4. Praise attributions	3.00	1.55	.09 [-.02, .19]	.19** [.09, .29]	.02 [-.08, .13]	
	5. Blame attributions	3.83	1.66	.37** [.28, .46]	.36** [.27, .45]	.08 [-.02, .19]	.19** [.09, .29]
Indeterministic universe (<i>n</i> = 355)	1. Free will attributions	5.36	0.87				
	2. Intent attributions	3.41	1.61	-.01 [-.12, .09]			
	3. Causality attributions	5.11	1.09	.36** [.26, .44]	.31** [.21, .40]		
	4. Praise attributions	3.17	1.64	-.03 [-.13, .08]	.20** [.09, .29]	.07 [-.03, .18]	
	5. Blame attributions	5.00	1.27	.24** [.14, .34]	.24** [.14, .34]	.22** [.11, .31]	.10 [-.00, .20]

Table S10

Study 1: Results of 3 x 2 ANOVA testing the effects of outcome and universe on attributions intentionality (combined dv, split dvs)

Intentionality attribution (two items combined)					Intentionality attribution (item 1)				Knowledge attribution (item 1)			
Factor	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p
Outcome (Help vs Harm)	182.80	1	<.001	.37	125.83	1	<.001	.29	154.89	1	<.001	.34
Universe	2.55	2	.080	.02	1.13	2	.324	.01	2.97	2	.053	.02
Outcome × Universe	1.62	2	.201	.01	0.47	2	.628	.00	2.35	2	.097	.02

Note. Outcome and Universe are between subject variables. *df* = degree of freedom, η^2_p = partial eta-squared.

Table S11

Study 2: Results of 3 x 2 ANOVA testing the effects of outcome and universe on attributions intentionality (combined dv, split dvs)

Intentionality attribution (two items combined)					Intentionality attribution (item 2)				Knowledge attribution (item 2)			
Factor	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p	<i>F</i>	<i>df</i>	<i>p</i>	η^2_p
Outcome (Help vs Harm)	182.80	1	<.001	.37	480.64	1	<.001	.31	647.58	1	<.001	.37
Universe	2.55	2	.080	.02	0.25	2	.778	.00	3.19	2	.041	.01
Outcome × Universe	1.62	2	.201	.01	1.50	2	.223	.00	1.80	2	.165	.00

Note. Outcome and Universe are between subject variables. *df* = degree of freedom, η^2p = partial eta-squared.

References

- Allaire J, Xie Y, McPherson J, Luraschi J, Ushey K, Atkins A, Wickham H, Cheng J, Chang W, Iannone R (2021). *rmarkdown: Dynamic Documents for R*. R package version 2.8, <https://github.com/rstudio/rmarkdown>.
- Auguie B (2017). *gridExtra: Miscellaneous Functions for "Grid" Graphics*. R package version 2.3, <https://CRAN.R-project.org/package=gridExtra>.
- Bates D, Mächler M, Bolker B, Walker S (2015). “Fitting Linear Mixed-Effects Models Using lme4.” *Journal of Statistical Software*, **67**(1), 1–48. doi: [10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01).
- Bates D, Maechler M (2021). *Matrix: Sparse and Dense Matrix Classes and Methods*. R package version 1.3-4, <https://CRAN.R-project.org/package=Matrix>.
- Chang W, Cheng J, Allaire J, Sievert C, Schloerke B, Xie Y, Allen J, McPherson J, Dipert A, Borges B (2021). *shiny: Web Application Framework for R*. R package version 1.6.0, <https://CRAN.R-project.org/package=shiny>.
- Comtois D (2021). *summarytools: Tools to Quickly and Neatly Summarize Data*. R package version 0.9.9, <https://CRAN.R-project.org/package=summarytools>.
- Dahl D, Scott D, Roosen C, Magnusson A, Swinton J (2019). *xtable: Export Tables to LaTeX or HTML*. R package version 1.8-4, <https://CRAN.R-project.org/package=xtable>.
- Demin G (2020). *expss: Tables, Labels and Some Useful Functions from Spreadsheets and 'SPSS' Statistics*. R package version 0.10.7, <https://CRAN.R-project.org/package=expss>.
- Falissard B (2012). *psy: Various procedures used in psychometry*. R package version 1.1, <https://CRAN.R-project.org/package=psy>.
- Firke S (2021). *janitor: Simple Tools for Examining and Cleaning Dirty Data*. R package version 2.1.0, <https://CRAN.R-project.org/package=janitor>.
- Fox J, Weisberg S (2019). *An R Companion to Applied Regression*, Third edition. Sage, Thousand Oaks CA. <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>.
- Fox J, Weisberg S, Price B (2020). *carData: Companion to Applied Regression Data Sets*. R package version 3.0-4, <https://CRAN.R-project.org/package=carData>.
- Fual, F., Erdfelder, E., Lang, A., & Buchner, A. (2007). G*Power: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, **39**(2), 175–191. <https://doi.org/10.3758/BF03193146>

- Genz A, Bretz F (2009). *Computation of Multivariate Normal and t Probabilities*, series Lecture Notes in Statistics. Springer-Verlag, Heidelberg. ISBN 978-3-642-01688-2.
- Genz A, Bretz F, Miwa T, Mi X, Leisch F, Scheipl F, Hothorn T (2021). *mvtnorm: Multivariate Normal and t Distributions*. R package version 1.1-2, <https://CRAN.R-project.org/package=mvtnorm>.
- Graves S, Piepho H, Dorai-Raj LSwfhS (2019). *multcompView: Visualizations of Paired Comparisons*. R package version 0.1-8, <https://CRAN.R-project.org/package=multcompView>.
- Gronau Q (2020). *abtest: Bayesian A/B Testing*. R package version 0.2.2, <https://CRAN.R-project.org/package=abtest>.
- Harrell Jr FE, Dupont wcfC, others. m (2021). *Hmisc: Harrell Miscellaneous*. R package version 4.5-0, <https://CRAN.R-project.org/package=Hmisc>.
- Henry L, Wickham H (2020). *purrr: Functional Programming Tools*. R package version 0.3.4, <https://CRAN.R-project.org/package=purrr>.
- Hervé M (2021). *RVAideMemoire: Testing and Plotting Procedures for Biostatistics*. R package version 0.9-79, <https://CRAN.R-project.org/package=RVAideMemoire>.
- Hlavac M (2018). *stargazer: Well-Formatted Regression and Summary Statistics Tables*. Central European Labour Studies Institute (CELSI), Bratislava, Slovakia. R package version 5.2.2, <https://CRAN.R-project.org/package=stargazer>.
- Hope RM (2013). *Rmisc: Rmisc: Ryan Miscellaneous*. R package version 1.5, <https://CRAN.R-project.org/package=Rmisc>.
- Hothorn T (2019). *TH.data: TH's Data Archive*. R package version 1.0-10, <https://CRAN.R-project.org/package=TH.data>.
- Hothorn T, Bretz F, Westfall P (2008). “Simultaneous Inference in General Parametric Models.” *Biometrical Journal*, **50**(3), 346–363.
- Kassambara A (2020). *ggpubr: 'ggplot2' Based Publication Ready Plots*. R package version 0.4.0, <https://CRAN.R-project.org/package=ggpubr>.
- Kassambara A (2021). *rstatix: Pipe-Friendly Framework for Basic Statistical Tests*. R package version 0.7.0, <https://CRAN.R-project.org/package=rstatix>.
- Kelley K (2020). *MBESS: The MBESS R Package*. R package version 4.8.0, <https://CRAN.R-project.org/package=MBESS>.

- Kuhn M, Jackson S, Cimentada J (2020). *corr: Correlations in R*. R package version 0.4.3, <https://CRAN.R-project.org/package=corr>.
- Lakens D (2017). “Equivalence tests: A practical primer for t-tests, correlations, and meta-analyses.” *Social Psychological and Personality Science*, **1**, 1–8. doi: [10.1177/1948550617697177](https://doi.org/10.1177/1948550617697177).
- Larmarange J (2021). *labelled: Manipulating Labelled Data*. R package version 2.8.0, <https://CRAN.R-project.org/package=labelled>.
- Lenth R (2021). *emmeans: Estimated Marginal Means, aka Least-Squares Means*. R package version 1.6.1, <https://CRAN.R-project.org/package=emmeans>.
- Lenth RV (2016). “Least-Squares Means: The R Package lsmeans.” *Journal of Statistical Software*, **69**(1), 1–33. doi: [10.18637/jss.v069.i01](https://doi.org/10.18637/jss.v069.i01).
- Long JA (2019). *interactions: Comprehensive, User-Friendly Toolkit for Probing Interactions*. R package version 1.1.0, <https://cran.r-project.org/package=interactions>.
- Long JA (2020). *jtools: Analysis and Presentation of Social Scientific Data*. R package version 2.1.0, <https://cran.r-project.org/package=jtools>.
- Lüdtke D (2019). *esc: Effect Size Computation for Meta Analysis (Version 0.5.1)*. doi: [10.5281/zenodo.1249218](https://doi.org/10.5281/zenodo.1249218), <https://CRAN.R-project.org/package=esc>.
- Lüdtke D (2021). *sjlabelled: Labelled Data Utility Functions (Version 1.1.8)*. doi: [10.5281/zenodo.1249215](https://doi.org/10.5281/zenodo.1249215), <https://CRAN.R-project.org/package=sjlabelled>.
- Lüdtke D, Ben-Shachar M, Patil I, Waggoner P, Makowski D (2021). “performance: An R Package for Assessment, Comparison and Testing of Statistical Models.” *Journal of Open Source Software*, **6**(60), 3139. doi: [10.21105/joss.03139](https://doi.org/10.21105/joss.03139).
- Mair P, Wilcox R (2020). “Robust Statistical Methods in R Using the WRS2 Package.” *Behavior Research Methods*, **52**, 464–488.
- Mangiafico S (2021). *rcompanion: Functions to Support Extension Education Program Evaluation*. R package version 2.4.1, <https://CRAN.R-project.org/package=rcompanion>.
- Morey R, Rouder J (2018). *BayesFactor: Computation of Bayes Factors for Common Designs*. R package version 0.9.12-4.2, <https://CRAN.R-project.org/package=BayesFactor>.
- Müller K, Wickham H (2021). *tibble: Simple Data Frames*. R package version 3.1.2, <https://CRAN.R-project.org/package=tibble>.

- Navarro D (2015). *Learning statistics with R: A tutorial for psychology students and other beginners. (Version 0.5)*. University of Adelaide, Adelaide, Australia. R package version 0.5, <http://ua.edu.au/ccs/teaching/lsr>.
- Ogle DH, Wheeler P, Dinno A (2021). *FSA: Fisheries Stock Analysis*. R package version 0.8.32, <https://github.com/droglenc/FSA>.
- Ooms J (2021). *magick: Advanced Graphics and Image-Processing in R*. R package version 2.7.2, <https://CRAN.R-project.org/package=magick>.
- Patil I (2021). “Visualizations with statistical details: The 'ggstatsplot' approach.” *PsyArxiv*. doi: [10.31234/osf.io/p7mku](https://doi.org/10.31234/osf.io/p7mku), <https://psyarxiv.com/p7mku/>.
- Plummer M, Best N, Cowles K, Vines K (2006). “CODA: Convergence Diagnosis and Output Analysis for MCMC.” *R News*, 6(1), 7–11. <https://journal.r-project.org/archive/>.
- R Core Team (2020). *foreign: Read Data Stored by 'Minitab', 'S', 'SAS', 'SPSS', 'Stata', 'Systat', 'Weka', 'dBase',* R package version 0.8-81, <https://CRAN.R-project.org/package=foreign>.
- Revelle W (2021). *psych: Procedures for Psychological, Psychometric, and Personality Research*. Northwestern University, Evanston, Illinois. R package version 2.1.3, <https://CRAN.R-project.org/package=psych>.
- Robinson D, Hayes A, Couch S (2021). *broom: Convert Statistical Objects into Tidy Tibbles*. R package version 0.7.6, <https://CRAN.R-project.org/package=broom>.
- Sarkar D (2008). *Lattice: Multivariate Data Visualization with R*. Springer, New York. ISBN 978-0-387-75968-5, <http://lmdvr.r-forge.r-project.org>.
- Selker R, Love J, Dropmann D (2020). *jmv: The 'jamovi' Analyses*. R package version 1.2.23, <https://CRAN.R-project.org/package=jmv>.
- Stanley D (2021). *apaTables: Create American Psychological Association (APA) Style Tables*. R package version 2.0.8, <https://CRAN.R-project.org/package=apaTables>.
- Terry M. Therneau, Patricia M. Grambsch (2000). *Modeling Survival Data: Extending the Cox Model*. Springer, New York. ISBN 0-387-98784-3.
- Therneau T (2021). *A Package for Survival Analysis in R*. R package version 3.2-11, <https://CRAN.R-project.org/package=survival>.
- Urbanek S (2013). *png: Read and write PNG images*. R package version 0.1-7, <https://CRAN.R-project.org/package=png>.
- Venables WN, Ripley BD (2002). *Modern Applied Statistics with S*, Fourth edition. Springer, New York. ISBN 0-387-95457-0, <https://www.stats.ox.ac.uk/pub/MASS4/>.

- Wickham H (2007). “Reshaping Data with the reshape Package.” *Journal of Statistical Software*, **21**(12), 1–20. <http://www.jstatsoft.org/v21/i12/>.
- Wickham H (2011). “The Split-Apply-Combine Strategy for Data Analysis.” *Journal of Statistical Software*, **40**(1), 1–29. <http://www.jstatsoft.org/v40/i01/>.
- Wickham H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. ISBN 978-3-319-24277-4, <https://ggplot2.tidyverse.org>.
- Wickham H (2019). *stringr: Simple, Consistent Wrappers for Common String Operations*. R package version 1.4.0, <https://CRAN.R-project.org/package=stringr>.
- Wickham H (2021). *forcats: Tools for Working with Categorical Variables (Factors)*. R package version 0.5.1, <https://CRAN.R-project.org/package=forcats>.
- Wickham H (2021). *tidyr: Tidy Messy Data*. R package version 1.1.3, <https://CRAN.R-project.org/package=tidyr>.
- Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Grolemund G, Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM, Müller K, Ooms J, Robinson D, Seidel DP, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H (2019). “Welcome to the tidyverse.” *Journal of Open Source Software*, **4**(43), 1686. doi: [10.21105/joss.01686](https://doi.org/10.21105/joss.01686).
- Wickham H, Bryan J (2021). *usethis: Automate Package and Project Setup*. R package version 2.0.1, <https://CRAN.R-project.org/package=usethis>.
- Wickham H, François R, Henry L, Müller K (2021). *dplyr: A Grammar of Data Manipulation*. R package version 1.0.6, <https://CRAN.R-project.org/package=dplyr>.
- Wickham H, Hester J (2020). *readr: Read Rectangular Text Data*. R package version 1.4.0, <https://CRAN.R-project.org/package=readr>.
- Wickham H, Hester J, Chang W (2021). *devtools: Tools to Make Developing R Packages Easier*. R package version 2.4.2, <https://CRAN.R-project.org/package=devtools>.
- Wickham H, Miller E (2021). *haven: Import and Export 'SPSS', 'Stata' and 'SAS' Files*. R package version 2.4.1, <https://CRAN.R-project.org/package=haven>.
- Xie Y (2014). “knitr: A Comprehensive Tool for Reproducible Research in R.” In Stodden V, Leisch F, Peng RD (eds.), *Implementing Reproducible Computational Research*. Chapman and Hall/CRC. ISBN 978-1466561595, <http://www.crcpress.com/product/isbn/9781466561595>.
- Xie Y (2015). *Dynamic Documents with R and knitr*, 2nd edition. Chapman and Hall/CRC, Boca Raton, Florida. ISBN 978-1498716963, <https://yihui.org/knitr/>.

- Xie Y (2016). *bookdown: Authoring Books and Technical Documents with R Markdown*. Chapman and Hall/CRC, Boca Raton, Florida. ISBN 978-1138700109, <https://bookdown.org/yihui/bookdown>.
- Xie Y (2021). *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.22, <https://github.com/rstudio/bookdown>.
- Xie Y (2021). *knitr: A General-Purpose Package for Dynamic Report Generation in R*. R package version 1.33, <https://yihui.org/knitr/>.
- Xie Y, Allaire J, Golemund G (2018). *R Markdown: The Definitive Guide*. Chapman and Hall/CRC, Boca Raton, Florida. ISBN 9781138359338, <https://bookdown.org/yihui/rmarkdown>.
- Xie Y, Dervieux C, Riederer E (2020). *R Markdown Cookbook*. Chapman and Hall/CRC, Boca Raton, Florida. ISBN 9780367563837, <https://bookdown.org/yihui/rmarkdown-cookbook>.
- Zeileis A, Croissant Y (2010). “Extended Model Formulas in R: Multiple Parts and Multiple Responses.” *Journal of Statistical Software*, **34**(1), 1–13. doi: [10.18637/jss.v034.i01](https://doi.org/10.18637/jss.v034.i01).
- Zeileis A, Grothendieck G (2005). “zoo: S3 Infrastructure for Regular and Irregular Time Series.” *Journal of Statistical Software*, **14**(6), 1–27. doi: [10.18637/jss.v014.i06](https://doi.org/10.18637/jss.v014.i06).
- Zeileis A, Hothorn T (2002). “Diagnostic Checking in Regression Relationships.” *R News*, **2**(3), 7–10. <https://CRAN.R-project.org/doc/Rnews/>.
- Zhu H (2021). *kableExtra: Construct Complex Table with 'kable' and Pipe Syntax*. R package version 1.3.4, <https://CRAN.R-project.org/package=kableExtra>.

Peer Review and Communication History

MS Title: Asymmetries in Attributions of Blame and Praise, Intent, and Causality: Free Will, Responsibility, and the Side-effect Effect

Author Names: Adrien Fillon, Subramanya Prasad Chandrashekar, Gilad Feldman

Submitted: Feb 5, 2023

Editor First Decision: Revise & Resubmit
Jun 6, 2023

Dear Dr. Feldman,

I have now received 3 reviews of your manuscript, “Asymmetries in attributions of blame and praise, intent, and causality: Free will, responsibility, and the side-effect effect”, from researchers with special expertise in moral psychology and philosophy. The reviewers had mixed reactions to your manuscript. I agree with the reviewers that your manuscript has important strengths and also that there are some issues that need to be addressed. I therefore encourage you to submit a revised version for further consideration at Collabra: Psychology.

The reviewers did an outstanding job in their reviews. I will highlight issues I think are particularly salient here. In your resubmission, please include a document with a point-by-point response to both the points I list here and the reviewers’ comments, outlining each change made in your manuscript or providing a suitable rebuttal.

Major revisions:

- Both reviewer 1 (R1) and reviewer 3 (R3) bring up issues with how causality is being measured, and suggest a distinction between knowledge and causality. R3 brings up a similar concern regarding how intent is being introduced and assessed. These are important points, and ones that would likely need to be addressed with the inclusion of more data that assesses these constructs more directly.
- All reviewers expressed frustration with the lack of setup for the hypotheses in the introduction, and note places where review of the extant literature is underdeveloped/lacking. R1 notes several papers that will be useful in mitigating this concern.
- Reviewer 2 (R2) notes that there are a number of results, but that these are not interpreted nor well-introduced. I agree, and also note that this is particularly the case for the regret variable. The inclusion of this variable was

not sufficiently motivated prior to Study 2, which was confusing as a reader (a concern also brought up by R1).

- At several points you bring up the relative novelty of examining both praise and blame across both helpful and harmful actions – however, I found the praise results to be downright confusing, and not addressed nor interpreted for the reader. In both Studies 1 and 2, the praise ratings were higher for the harmful than helpful actions in the deterministic/indeterministic universes, but this was reversed for the control condition. While some of these differences are not statistically significant, based on the means, standard deviations, and sample sizes reported on pg. 27, some of them almost certainly are. To me, this potentially suggests that participants may have been more generally confused about the deterministic vs. indeterministic universe manipulation. Some way of assessing comprehension of this manipulation would be useful in mitigating this concern.

In summary, I think this is a promising manuscript and, I hope you will revise it for further consideration at Collabra: Psychology. I look forward to receiving your revision. Please see the instructions below for submitting your revision.

Please ensure that your revised files adhere to our author guidelines, and that the files are fully copyedited/proofed prior to upload. Please also ensure that all copyright permissions have been obtained. This may be the last opportunity for major editing, therefore please fully check your file prior to re-submission.

If you have any questions or difficulties during this process, please contact the editorial office at editorialoffice@collabra.org.

We hope you can submit your revision within the next six weeks. If you cannot make this deadline, please let us know as early as possible.

Sincerely,

Chelsea Helion

Reviewer 1

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs)				✓	

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
the authors wish to study). (Choose “Neutral” if this is not an empirical manuscript)					
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose “Neutral” if this is not an empirical manuscript)					✓
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose “Neutral” if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)			✓		

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

In two studies, the authors examined the side-effect effect (SEE) and how beliefs about determinism and free-will more strongly impact attributions of blame vs praise.

Most crucially, I found the Introduction to be relatively sparse and underspecified. The authors review some of the past work on the SEE, but it's unclear what open questions remain – and more importantly, why those open questions are worth investigating. What new insight would we gain about SEE and social/moral cognition more broadly through these studies? As is, the Introduction feels atheoretical and it's unclear what specific predictions and hypotheses the authors have. Perhaps the closest is on pp. 4-5, “we speculated that an agent in a situation involving a harmful outcome scenario, even when the outcome was a side-effect, is attributed more free will than an agent with a beneficial outcome.” This is a good start, but I believe that more clarity regarding why the authors investigated these questions, let alone what they think the potential answers are, would do much to improve the paper. Also, the previous sentence seems to suggest that the general question being posed by this research has already been answered: “Empirical studies found support for the view that actions and outcomes of negative valence led to higher attribution of free will to the agent than outcomes of positive valence, even for non-moral scenarios (Feldman et al., 2016; Fillon et al., 2021)” (p. 4). Given that past research has already found that people attribute higher free will for

negative outcomes than positive outcomes, what does the present research add above and beyond that past research?

As part of this concern about the Introduction, a couple of theory papers seem especially relevant here. Malle et al., 2014 provide a lengthy review and theory for moral blame, while Anderson et al., 2020 offer both a theory of moral praise and of when/why praise-blame asymmetries should occur. Integrating these theory papers seems like it would help clarify how people make judgments of blame and praise and therefore why asymmetries should occur.

Likewise, the authors emphasize that they are bringing together two well-known experimental philosophy paradigms, but it's unclear why they are doing so. Again, what do we hope to learn from studying SEE and free-will attributions in the same set of studies? What do beliefs about determinism tell us about this association?

Study 1: Causality attributions – Is this really the best term? The measures ask about knowing and understanding, which are not the same thing as being causally responsible. For example, I can know and understand the potential effects of the new program on the environment, but since I'm not the chairman, I don't think people would hold me causally responsible. That is, these items seem to be getting at knowledge and not causal responsibility. Why not just ask a more face-valid question, like "In Universe D, the chairman caused the effects of the new program on the environment"?

Study 2 – Why measure regret? The only previous mention of regret is on p. 5 noting a correspondence between regret and free will attributions. What would regret tell us?

How do the authors interpret their finding that attributions of praise were similar when the outcome was both positive and negative? That is, why would a harmful act be as praiseworthy as a positive act? This seems to suggest that participants are not treating the measure of praise in a moral sense, but perhaps in a performance sense (i.e., good for the company even if harmful to the environment and therefore laudable). If that's the case, it's harder to make sense of the praise results.

I have similar comments about the GD as I did about the Introduction – it's hard to tell what exactly the authors are arguing for and what we know now about attributions of blame/praise, intentionality, and free-will that we didn't know before. The authors point to a number of existing theories/hypotheses and say that their results confirm those theories/hypotheses (eg "In conclusion, our results confirmed the strong relationship between attribution of blame and free will.", p. 45). This is helpful, but it doesn't suggest that anything new has been learned. I'm not saying that all papers need to be super novel to justify publication, but I'm left

wondering why I should read (let alone cite) this set of studies that past research doesn't already tell us.

Another branch of research that the authors could consider citing and engaging with is the work on moral character, which argues that people make judgments not just of acts but of the people involved in those acts (see work by Geoff Goodwin, David Pizarro, Kurt Gray, and others). Such work offers reinterpretations and explanations for the some of the positive-negative asymmetry – for example, attributions to “true selves” (eg Newman et al., 2015) and assumptions about underlying motivations and desires.

Minor comment:

p. 2 – In the SSE vignette, the second paragraph should have a bracket for “I don't care at all about [positive condition: helping; negative condition: harming] the environment.”

References:

Anderson, R. A., Crockett, M. J., & Pizarro, D. A. (2020). A theory of moral praise.

Trends in Cognitive Sciences, 24(9), 694–

703. <https://doi.org/10.1016/j.tics.2020.06.008>

Malle, B. F., Guglielmo, S., & Monroe, A. E. (2014). A theory of blame. Psychological

Inquiry, 25(2), 147–186. <https://doi.org/10.1080/1047840X.2014.877340>

Newman, G. E., De Freitas, J., & Knobe, J. (2015). Beliefs About the True Self Explain Asymmetries Based on Moral Judgment. Cognitive Science, 39(1), 96–

125. <https://doi.org/10.1111/cogs.12134>

Reviewer 2

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose “Neutral” if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose “Neutral” if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose “Neutral” if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

This is an interesting paper examining possible connections between the Side Effect Effect (SEE) and free will attributions. The two studies essentially take an “exploring small, confirming big” approach (Sakaluk, 2016), which I find compelling (and underutilized). Overall, the results are sensible and consistent with past findings, if not especially novel. I have some suggestions for how to improve the paper:

-It would be very helpful if the authors would formally, explicitly state their hypotheses in the Introduction section. Right now, predictions are alluded to and vaguely mentioned, but they are never clearly and precisely stated. This makes it difficult to follow the (many) analyses and what they are meant to test.

-I am a bit concerned about the Free Will measure. The items ask whether the chairman “had to” do as he did, or was “free to” do otherwise. The potential problem is that this could plausibly tap other constraints on behavior, aside from free will per se. To give one example, one could plausibly argue that the chairman has a fiduciary duty to maximize returns for his shareholders, and so he “had to” implement a program that would increase profits, and he was therefore not “free to” do otherwise. Of course, this measure does show the predicted difference across the manipulation of universe type, but that doesn’t mean that it is necessarily only measuring what it is meant to. This issue of discriminant validity should at least be mentioned and discussed.

-There are many, many analyses, and little interpretation of the results. This makes it very hard to follow what is going on in the Results sections and make sense of what it all means. Some signposting would be very helpful. To give one particular example, on p. 36, the authors mention a significant three-way interaction between universe, outcome, and measure, but never elaborate on the pattern of results or interpret what this interaction means.

-I think it is too strong to claim that “we found support for SEE generalizability to free will” (p. 42). The classic SEE was only replicated for the free will measure in Study 2, not Study 1, and the effect size is quite small. Similar statements are made elsewhere. The authors should be careful about over-interpreting/going beyond the data.

Reviewer 3

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose “Neutral” if this is not an empirical manuscript)	✓				
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose “Neutral” if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose “Neutral” if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

As a philosopher, the thing I found most frustrating about this paper is that it wasn't make clear enough at the beginning what specific hypotheses the experiments were designed to test. That is, what are the different views about free will, determinism, intentionality, and causality, and about how these concepts are connected in our thinking, that generate various hypotheses that the authors' experiments could then help test? It's possible to piece some of this together, especially from what's said at the end of the paper, but nowhere was it laid out very clearly, and because of this it felt like the narrative of the paper was too much along the lines of “We wanted to see if these things would happen, so we ran some experiments ...” As a result, the results of the studies came across mainly as a string of numbers.

I also have some worries about the design of the studies themselves, specifically concerning the intentionality and causality attributions.

First, for the intentionality attributions the first item presented for agreement or disagreement (“In Universe D, the chairman intentions were to have such implications of the new program on the environment?”) is simply ungrammatical, so I don’t see what can be learned from participants’ responses to it. Further, it isn’t clear (at least from what I saw) whether here and elsewhere the authors then combined the responses to these two items into a single variable – though if they did then that would be a strange decision, given that the first item concerns the chairman’s /intentions/ while the second is rather about what he /intentionally did/, and from the perspective of philosophy these are quite different matters.

Second, the situation with the causality attributions is even worse since the two items presented have to do only with the chairman’s /knowledge/ and /understanding/, rather than with what he caused to happen. While there is indeed some evidence that knowledge attributions are also subject to a norm effect, these items do nothing at all to measure attributions of causality, which is what the authors aimed to study.

Author Response

Jun 19, 2024

Reply to Collabra decision letter: **Free will and the side-effect effect**

We would like to thank the editor and the reviewers for their useful suggestions and below we provide a detailed response as well as a tally of all the changes that were made in the manuscript. For an easier overview of all the changes made, we also provide a summary of changes. Please note that the editor’s and reviewers’ comments are in bold with our reply underneath in normal script.

A track-changes comparison of the previous submission and the revised submission can be found on: <https://draftable.com/compare/WqOKbIArgczu>

Summary of changes

Below we provide a table with a summary of the main changes to the manuscript and our response to the editor and reviewers:

Section	Actions taken in the current manuscript
General	All: modification of causality label to knowledge
Introduction	<p>All: Added a chunk regarding motivation to test regret. Added a new table for the hypotheses.</p> <p>R1, R3: Addition of a new section “linking free will and the side effect effect”.</p> <p>R2: Addition of a paragraph focusing on praise.</p>
Methods	Ed, R3: added a footnote regarding intentionality.
Results	All: Added a reference to the hypotheses in the table to improve clarity.
Discussion	<p>All: Added a table summarizing the results obtained.</p> <p>Ed, R3: Added a limitation part explaining our view on intentionality.</p> <p>R1: Added a chunk regarding praise.</p> <p>R2: added a small chunk regarding motivation and free will.</p>
<i>Note.</i> Ed = Editor, R1/R2/R3 = Reviewer 1/2/3	

Reply to Editor: Dr./Prof. Chelsea Helion

I have now received 3 reviews of your manuscript, “Asymmetries in attributions of blame and praise, intent, and causality: Free will, responsibility, and the side-effect effect”, from researchers with special expertise in moral psychology and philosophy. The reviewers had mixed reactions to your manuscript. I agree with the reviewers that your manuscript has important strengths and also that there are some issues that need to be addressed. I therefore encourage you to submit a revised version for further consideration at Collabra: Psychology.

The reviewers did an outstanding job in their reviews. I will highlight issues I think are particularly salient here. In your resubmission, please include a document with a point-by-point response to both the points I list here and the reviewers’ comments, outlining each change made in your manuscript or providing a suitable rebuttal.

Thank you for the reviews obtained, your feedback, and the invitation to revise and resubmit.

Major revisions:

.1. Both reviewer 1 (R1) and reviewer 3 (R3) bring up issues with how causality is being measured, and suggest a distinction between knowledge and causality. R3 brings up a similar concern regarding how intent is being introduced and assessed. These are important points, and ones that would likely need to be addressed with the inclusion of more data that assesses these constructs more directly.

We agree and appreciate the feedback. Thank you and the reviewers for helping us realize we should have done better in explaining where these measures came from and to better articulate their meaning, context, and use.

The two questions were adapted from a previous published article: Beebe and Jensen (2012) in which the term causality is not mentioned, only the term knowledge. In the previous version of the MS, we referred to it as a close approximation of causality. We made changes throughout to refer to knowledge rather than to causality.

For the intention measure, we now explain it in detail in the manuscript, in the method section and in a newly introduced “limitations” section. We note that we conducted a series of analyses leading us to think that the grammatical error for the first item of intention does not mean the item is useless and that we think it can still be aggregated with the second item to analyze intention. However, as noted in the limitation section, one can still expand the knowledge of intention with a better, more precise question than ours.

Our understanding of the reviewers’ points was that none of them suggested additional data collection, and we decided not to conduct additional data collection. We believe that these issues should not have any crucial impact on the findings or the package overall, and that these points are best addressed with greater transparency and a discussion of the identified limitations.

References:

- Beebe, J. R., & Jensen, M. (2012). Surprising connections between knowledge and action: The robustness of the epistemic side-effect effect. *Philosophical psychology*, 25(5), 689-715.

.2. All reviewers expressed frustration with the lack of setup for the hypotheses in the introduction, and note places where review of the extant literature is underdeveloped/lacking. R1 notes several papers that will be useful in mitigating this concern.

We understand and appreciate the feedback urging us to do better.

We added a table in the introduction explaining the hypotheses, their rationale, with an indication of whether the hypotheses were exploratory or confirmatory (pre-registered), and the effect size for results in Studies 1 and 2. We also completely reworked the introduction, adding several subtitles regarding the SEE effect, the relationship between SEE and Free Will, and a dedicated part on the extension for regret.

We also completely reworked the discussion section, including several subsections to discuss every aspect of our findings; the side effect effect replication, the relationship with free will and the aspect of accountability, the extension on regret, and the new finding of the SEE regardless of the outcome.

Finally, we brought in the points from R1's suggested papers. We accommodated these suggestions both in the introduction and discussion sections of the manuscript. We believe that the revised manuscript is much improved, and readers should be able to easily follow the theoretical arguments. We are grateful for all the constructive positive suggestions.

3. Reviewer 2 (R2) notes that there are a number of results, but that these are not interpreted nor well-introduced. I agree, and also note that this is particularly the case for the regret variable. The inclusion of this variable was not sufficiently motivated prior to Study 2, which was confusing as a reader (a concern also brought up by R1).

We created an independent section on regret in the introduction and in the discussion to cover the regret variable. We modified the results section to improve readability, reducing the overall number of figures from 18 to 10. We also make the reading of the results consistent across variables and studies. Finally, we modified the method section, especially regarding our method to measure Praise and Blame to explain how it differs from the original method.

.4. At several points you bring up the relative novelty of examining both praise and blame across both helpful and harmful actions – however, I found the praise results to be downright confusing, and not addressed nor interpreted for the reader. In both Studies 1 and 2, the praise ratings were higher for the harmful than helpful actions in the deterministic/indeterministic universes, but this was reversed for the control condition. While some of these differences are not statistically significant, based on the means, standard deviations, and sample

sizes reported on pg. 27, some of them almost certainly are. To me, this potentially suggests that participants may have been more generally confused about the deterministic vs. indeterministic universe manipulation. Some way of assessing comprehension of this manipulation would be useful in mitigating this concern.

Thank you very much for the feedback. We take this remark as a complement to the remark by Reviewer 2 on Anderson's (2020) study on praise.

We believe that our reporting may have caused confusion regarding what attributions were contrasted, given that we collected both praise for a possible positive outcome and blame for a possible negative outcome (within-subject design) regardless of the outcome condition that the participant was assigned to. In the previous version we simplified the reporting to a 3x2 ignoring the within-subject measures factor, and only examining praise for positive and blame for negative, but overlooking the within and not reporting the fuller 3x2x2 may result in confusion. Therefore, the revised manuscript addresses that and now reports and plots the full 3x2x2. We now make the blame and praise measures and the interactions clearer with full reporting of both blame and praise across all conditions.

In summary, I think this is a promising manuscript and, I hope you will revise it for further consideration at Collabra: Psychology. I look forward to receiving your revision. Please see the instructions below for submitting your revision.

Thank you very much for all the feedback and help to improve the manuscript.

Reply to Reviewer #1

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose “Neutral” if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose “Neutral” if this is not an empirical manuscript)					✓
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose “Neutral” if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)			✓		

In two studies, the authors examined the side-effect effect (SEE) and how beliefs about determinism and free-will more strongly impact attributions of blame vs praise.

.1. Most crucially, I found the Introduction to be relatively sparse and underspecified. The authors review some of the past work on the SEE, but it’s unclear what open questions remain – and more importantly, why those open questions are worth investigating. What new insight would we gain about SEE and social/moral cognition more broadly through these studies? As is, the Introduction feels atheoretical and it’s unclear what specific predictions and hypotheses the authors have. Perhaps the closest is on pp. 4-5, “we speculated that an agent in a situation involving a harmful outcome scenario, even when the outcome was a side-effect, is attributed more free will than an agent with a beneficial outcome.” This is a good start, but I believe that more clarity

regarding why the authors investigated these questions, let alone what they think the potential answers are, would do much to improve the paper. Also, the previous sentence seems to suggest that the general question being posed by this research has already been answered: “Empirical studies found support for the view that actions and outcomes of negative valence led to higher attribution of free will to the agent than outcomes of positive valence, even for non-moral scenarios (Feldman et al., 2016; Fillon et al., 2021)” (p. 4). Given that past research has already found that people attribute higher free will for negative outcomes than positive outcomes, what does the present research add above and beyond that past research?

Thank you very much for the feedback. We considerably reworked the document, especially the introduction part to explain the three main contributions of this manuscript: the replication of the side-effect effect, the relationships with free-will and the extension showing that the side-effect is not dependent on the positivity or negativity of the direct effect. Every development has its own section dedicated to explaining the phenomenon studied.

We also added a summary Table 1 for the hypotheses, explaining them, and adding a rationale to them. Finally, we clearly categorized all analyses, labeling those as either confirmatory or exploratory to make the contributions clearer.

For the last quoted sentence, we added that the relationship between free will and negative outcomes was not shown with the side-effect effect experimentation. Using the side-effect effect scenario can increase the scope of the theory by indicating robustness of the relationship.

We hope that these changes help the reader understand better our motivation, how we planned to analyze the results and what the important results are for the theory.

As part of this concern about the Introduction, a couple of theory papers seem especially relevant here. Malle et al., 2014 provide a lengthy review and theory for moral blame, while Anderson et al., 2020 offer both a theory of moral praise and of when/why praise-blame asymmetries should occur. Integrating these theory papers seems like it would help clarify how people make judgments of blame and praise and therefore why asymmetries should occur.

Thank you for these references. We added them in the introduction section. Anderson's paper was especially helpful, and based on that, we added a short section in the discussion regarding the differences between praise and blame.

Likewise, the authors emphasize that they are bringing together two well-known experimental philosophy paradigms, but it's unclear why they are doing so. Again, what do we hope to learn from studying SEE and free-will attributions in the same set of studies? What do beliefs about determinism tell us about this association?

We appreciate that feedback. In our revision, we have worked to elaborate on the link between the two, examining the links between intent, free will, and accountability. The key principle relationship here is accountability: Do we make the actor more accountable because he was free to do otherwise? We added a full section explaining this mechanism in the introduction and in the discussion.

Study 1: Causality attributions – Is this really the best term? The measures ask about knowing and understanding, which are not the same thing as being causally responsible. For example, I can know and understand the potential effects of the new program on the environment, but since I'm not the chairman, I don't think people would hold me causally responsible. That is, these items seem to be getting at knowledge and not causal responsibility. Why not just ask a more face-valid question, like "In Universe D, the chairman caused the effects of the new program on the environment"?

Thank you for this valuable feedback.

The label "causality" was used, as a replication of the omission bias referring to similar items as "causality" (Jamison et al., 2020). We believe that this measure was adapted from Beebe and Jensen (2012) who indeed referred to "knowledge" rather than "causality". We agree that "knowledge" is a more accurate label and have adjusted it throughout the manuscript.

We now also more clearly refer to Beebe and Jensen's work and explain our findings in the context of their article, providing confirmation of their findings with an added note in the discussion section.

References:

- Beebe, J. R., & Jensen, M. (2012). Surprising connections between knowledge and action: The robustness of the epistemic side-effect effect. *Philosophical psychology*, 25(5), 689-715.

- Jamison, J., Yay, T., & Feldman, G. (2020). Action-inaction asymmetries in moral scenarios: Replication of the omission bias examining morality and blame with extensions linking to causality, intent, and regret. *Journal of Experimental Social Psychology*, 89, 103977.

Study 2 – Why measure regret? The only previous mention of regret is on p. 5 noting a correspondence between regret and free will attributions. What would regret tell us?

Thank you for this feedback. Your comments have helped us realize that readers would find the current explanation of regret incomplete. We added a dedicated section in the introduction explaining the regret extension “Extension: Attributions of regret and moral responsibility”, and a dedicated section in discussion in which we explain and interpret our results regarding regret with more details.

How do the authors interpret their finding that attributions of praise were similar when the outcome was both positive and negative? That is, why would a harmful act be as praiseworthy as a positive act? This seems to suggest that participants are not treating the measure of praise in a moral sense, but perhaps in a performance sense (i.e., good for the company even if harmful to the environment and therefore laudable). If that’s the case, it’s harder to make sense of the praise results.

Thank you. Your comment helped us realize that we have not sufficiently explained our design and the exact question asked. Actually, the attribution of praise was similar for a positive side effect of both a negative and positive main effect, so the side effect was always positive. The exact question is “In Universe D, the chairman should be applauded for his actions if they led to positive outcomes.” To be completely clear, we have not asked a question regarding praise for a negative action.

We decided to rework the paper entirely to ensure a clear comprehension of our operationalization of the side-effect. For example, we wrote in the result section of the original draft:

“We found support for a main effect of praise and blame, as blame attributions were higher than praise attributions.”

In the revised draft:

“Praise attributed for a positive side effect in the harmful outcome condition ($M = 2.91$, $SD = 1.76$) was similar to the praise attributed for a positive side effect in the helpful outcome condition ($M = 2.99$, $SD = 1.38$), yet blame attributed to the negative side effect in the harmful outcome condition was higher ($M = 5.02$, $SD = 1.32$) than in the helpful outcome condition ($M = 4.20$, $SD = 1.53$).”

I have similar comments about the GD as I did about the Introduction – it’s hard to tell what exactly the authors are arguing for and what we know now about attributions of blame/praise, intentionality, and free-will that we didn’t know before. The authors point to a number of existing theories/hypotheses and say

that their results confirm those theories/hypotheses (eg “In conclusion, our results confirmed the strong relationship between attribution of blame and free will.”, p. 45). This is helpful, but it doesn’t suggest that anything new has been learned. I’m not saying that all papers need to be super novel to justify publication, but I’m left wondering why I should read (let alone cite) this set of studies that past research doesn’t already tell us.

We appreciate this feedback. We extensively reworked the discussion section to 1) explain in more detail our results and 2) mirror the introduction regarding the sections and the hypotheses. We also created a table summarizing all our hypotheses in the introduction and the findings across the two studies (Table 1). We hope that the new version is clearer, easier to read, and provides more insightful explanations for the theories.

Another branch of research that the authors could consider citing and engaging with is the work on moral character, which argues that people make judgments not just of acts but of the people involved in those acts (see work by Geoff Goodwin, David Pizarro, Kurt Gray, and others). Such work offers reinterpretations and explanations for the some of the positive-negative asymmetry – for example, attributions to “true selves” (eg Newman et al., 2015) and assumptions about underlying motivations and desires.

Thank you very much for the suggestions. We implemented the previous citations but refrained from citing the work on the true selves as we did not see how it could enter the scope of our study. Our focus is on the relationship between perception of free will, accountability, and blaming for an action, so we are concerned that this is too far related to our main argument.

Minor comment:

p. 2 – In the SSE vignette, the second paragraph should have a bracket for “I don’t care at all about [positive condition: helping; negative condition: harming] the environment.”

Thank you for catching that, appreciated. We fixed it.

References:

Anderson, R. A., Crockett, M. J., & Pizarro, D. A. (2020). A theory of moral praise. *Trends in Cognitive Sciences*, 24(9), 694–703.

<https://doi.org/10.1016/j.tics.2020.06.008>

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

<https://doi.org/10.1080/1047840X.2014.877340>

Newman, G. E., De Freitas, J., & Knobe, J. (2015). Beliefs About the True Self Explain Asymmetries Based on Moral Judgment. *Cognitive Science*, 39(1), 96–125. <https://doi.org/10.1111/cogs.12134>

Reply to Reviewer #2

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)			✓		

This is an interesting paper examining possible connections between the Side Effect Effect (SEE) and free will attributions. The two studies essentially take an “exploring small, confirming big” approach (Sakaluk, 2016), which I find compelling (and underutilized). Overall, the results are sensible and consistent with past findings, if not especially novel.

Thank you for the supportive positive opening comment.

I have some suggestions for how to improve the paper:

-It would be very helpful if the authors would formally, explicitly state their hypotheses in the Introduction section. Right now, predictions are alluded to and vaguely mentioned, but they are never clearly and precisely stated. This makes it difficult to follow the (many) analyses and what they are meant to test.

Thank you, we agree. We created a table in the introduction section with all the hypotheses, their rationale, categorizing whether they are exploratory or confirmatory, and summarizing our findings in the two studies. This could help a future meta-analysis or systematic review on the subject.

-I am a bit concerned about the Free Will measure. The items ask whether the chairman “had to” do as he did, or was “free to” do otherwise. The potential problem is that this could plausibly tap other constraints on behavior, aside from free will per se. To give one example, one could plausibly argue that the chairman has a fiduciary duty to maximize returns for his shareholders, and so he “had to” implement a program that would increase profits, and he was therefore not “free to” do otherwise. Of course, this measure does show the predicted difference across the manipulation of universe type, but that doesn’t mean that it is necessarily only measuring what it is meant to. This issue of discriminant validity should at least be mentioned and discussed.

We appreciate you raising this concern. Please find below some arguments to explain our position.

First, the measures we used are the ones used largely in the field by studies regarding attributions of free will (such as in Feldman et al, 2016; Feldman and Chandrashekar, 2018). Empirically, the free will universe manipulation – which is specifically about the broader aspect of determinism - impacted these free will attributions. It means that there is a clear and direct link between the measure of free will attributions and the philosophical notions of free will.

Addressing the more conceptual issue you raised - per definition, free will is about the capacity for choice, both regarding internal and external constraints (Feldman, 2018). We agree that it would be interesting to disentangle internal and external constraints in this situation and to compare the contextual factors to the broad fundamental philosophical factors. These are exciting directions for future research in follow-up studies. We added it to our limitations for the generalization of the side-effect effect to free will:

“One limitation lies in our manipulation of the free will universe, as it refers very broadly to the ability to choose without disentangling the constraints underlying the inability to choose. In the discussion regarding free will, there are context-dependent constraints (e.g., job role) and broader, more fundamental factors that restrict choice that are close to the philosophical debate on free will (e.g., determinism, higher power, genes, physics, etc.). While the free will universe manipulation is close to the philosophical debate and the manipulation impacted free will attributions, it is possible that free will attributions might also be related to the contextual aspects of choice. Finally, there is the possibility that the universe scenarios do not work as intended, as participants can have difficulties understanding the consequences of a deterministic universe. Future studies can expand on our findings to examine more specific constraints and how the effects we reported vary depending on the type of constraint or operationalization of the universe.”

-There are many, many analyses, and little interpretation of the results. This makes it very hard to follow what is going on in the Results sections and make sense of what it all means. Some signposting would be very helpful. To give one particular example, on p. 36, the authors mention a significant three-way interaction between universe, outcome, and measure, but never elaborate on the pattern of results or interpret what this interaction means.

Thank you for this feedback. We have revised to address this issue by restructuring the results section and adding better signposting with summary tables in the introduction.

We completely reworked the introduction, method, result and discussion sections regarding the three-way-interaction to explain better what the core argument is related to this particular analysis. This is that the blame for the side-effect is stronger than praise, no matter the main effect (positive or negative). This result is a major contribution to the literature, for which previous studies only asked participants for either negative or positive side effects, and not for both.

An example of our rework regarding page 36 is:

we wrote in the result section of the original draft:

“We found support for a main effect of praise and blame, as blame attributions were higher than praise attributions.”

In the revised draft:

“Praise attributed for a positive side effect in the harmful outcome condition ($M = 2.91$, $SD = 1.76$) was similar to the praise attributed for a positive side effect in the helpful outcome condition ($M = 2.99$, $SD = 1.38$), yet blame attributed to the negative side effect in the harmful outcome condition was higher ($M = 5.02$, $SD = 1.32$) than in the helpful outcome condition ($M = 4.20$, $SD = 1.53$).”

-I think it is too strong to claim that “we found support for SEE generalizability to free will” (p. 42). The classic SEE was only replicated for the free will measure in Study 2, not Study 1, and the effect size is quite small. Similar statements are made elsewhere. The authors should be careful about over-interpreting/going beyond the data.

Thank you, we took care to modify this part and the rest of the discussion to not over-interpret our results. Our main findings are the strong correlation between blame and free-will attribution and the difference in the praise/blame in the indeterministic and deterministic universes.

Reply to Reviewer #3

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)	✓				
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

As a philosopher, the thing I found most frustrating about this paper is that it wasn't make clear enough at the beginning what specific hypotheses the experiments were designed to test. That is, what are the different views about free will, determinism, intentionality, and causality, and about how these concepts are connected in our thinking, that generate various hypotheses that the authors' experiments could then help test? It's possible to piece some of this together, especially from what's said at the end of the paper, but nowhere was it laid out very clearly, and because of this it felt like the narrative of the paper was too much along the lines of "We wanted to see if these things would happen, so we ran some experiments ..." As a result, the results of the studies came across mainly as a string of numbers.

Thank you for challenging us to do better. We heavily modified the manuscript, especially to increase the writing on the theories underlying our hypotheses in the introduction and interpret

them in the discussion. We created dedicated sections in the introduction for each specific hypothesis, and we summarized them (and the results) in Table 1. We hope that we now have sufficient arguments to support our analysis.

I also have some worries about the design of the studies themselves, specifically concerning the intentionality and causality attributions.

First, for the intentionality attributions the first item presented for agreement or disagreement (“In Universe D, the chairman intentions were to have such implications of the new program on the environment?”) is simply ungrammatical, so I don’t see what can be learned from participants’ responses to it.

Thank you very much for having caught it. We added a footnote in the method section and a limitation section explaining this problem as follows:

“We also found an oversight in the two items measuring attribution of intention, which were not grammatically clear. We adapted the two items on the intention from Jamison and colleagues’ (2020) study, which are not standardized and may have impacted one of the questions. However, the reliability coefficients were high for both studies. We conducted a 2x2 ANOVA on both items and noted that, even if the second item is higher than the first, they are affected the same way by the SEE and the type of universe. We added the results to OSF. We believe that despite the awkward phrasing, the two questions were similarly understood by the participants.”

Further, it isn’t clear (at least from what I saw) whether here and elsewhere the authors then combined the responses to these two items into a single variable – though if they did then that would be a strange decision, given that the first item concerns the chairman’s /intentions/ while the second is rather about what he /intentionally did/, and from the perspective of philosophy these are quite different matters.

Thank you. First, we note that for both studies, the reliability coefficient is high, indicating that the two items can be merged (we merged them by taking the mean of the two items). Second, based on your feedback, we conducted additional ANOVAs for both items and both studies. The results of this analysis can be found in Table S10 and S11 of the supplementary materials. We found a main effect of SEE on intention, no interaction effect with the universe, and no main effect of the universe except for intent2 in study 2, very close to the threshold of 0.05. The similarities between the two items in terms of results lead us to think that, even with the grammatical error, participants broadly understood that the two questions were similar. However, we noted this as a limitation and called for an improvement of this measure in the limitation section.

Second, the situation with the causality attributions is even worse since the two items presented have to do only with the chairman’s /knowledge/ and

/understanding/, rather than with what he caused to happen. While there is indeed some evidence that knowledge attributions are also subject to a norm effect, these items do nothing at all to measure attributions of causality, which is what the authors aimed to study.

Thank you. It is an oversight in our preregistration which then made its way to the main manuscript. We went back to the literature, and the questions regarding causality come from Beebe and Jensen (2012) in which the term “causality” is not used and the term knowledge is used instead. We modified the whole text regarding this term to use knowledge instead. We also modified the discussion section regarding the generalizability of the findings regarding knowledge.

Editor Second Decision: Revise & Resubmit
Sep 24, 2024

Dear Gilad Feldman,

I have now read your revised manuscript. I appreciate your careful attention to the concerns the reviewers and I raised. I am happy to provisionally accept your manuscript for submission. However, the reviewers found a few small things I would like you to address.

Minor revisions:

- Reviewer 2 recommends including details about how subjects were recruited (e.g., Mechanical Turk).
- Reviewers 2 & 3 noted the helpfulness of Table 1 (I’m noting it as well – it’s very clear and significantly enhances readability/comprehension of results!). However, Reviewer 3 (originally Reviewer 1) notes that the Introduction would still benefit from a brief discussion of the theoretical implications that motivated including the side effect effect and free will within the same experimental paradigm. I think this could possibly be well situated in the “Linking Free Will and SEE: Manipulation of both Agency and Outcome Valence” section, as currently that section explains how you manipulate the variables and that it hasn’t been done yet, but not what the theoretical implications would be were you to find what you predict (or its opposite, or a null). Alternatively, another place this could work well would be where you outline the order of studies under “The Present Investigation” – this section clearly outlines what you did, but not its potential implications for our understanding of either or both phenomena.

- It's a bit confusing on p. 15 when it reads "The deterministic and indeterministic universe conditions were presented as follows:" and then has a description of the control condition – I think this may be an error, and you instead want to reference Table 2?
- On page 40, you note that the correlation ($r = .36$) between free will and blame attributions is strong, but to my understanding, .36 would be indicative of a low-to-moderate correlation – I think you might want to state that it is significant instead?

I look forward to receiving your final revision and accepting it for publication in Collabra: Psychology.

Please ensure that your revised files adhere to our author guidelines, and that the files are fully copyedited/proofed prior to upload. Please also ensure that all copyright permissions have been obtained. This is the last opportunity for major editing, therefore please fully check your file prior to re-submission.

If you have any questions or difficulties during this process, please contact the editorial office at editorialoffice@collabra.org.

We hope you can submit your revision within the next six weeks. If you cannot make this deadline, please let us know as early as possible.

Sincerely,

Chelsea Helion

Reviewer 1

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

...

Reviewer 2

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

The authors have done a fine job of responding to my prior concerns. My worry about the discriminant validity of the free will measure remains, but at least the authors acknowledge it in the General Discussion now. Table 1 is exceptionally helpful in summarizing the research and guiding the reader.

There's only one issue that I noticed. I must have missed this last time, but Study 1 does not contain an explanation of where and how participants were recruited. I'm guessing (based on the reported demographics) that they were recruited via Mechanical Turk, as in Study 2, but this isn't stated explicitly. The authors should add this information.

Reviewer 3

Rating scale questions

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)					✓
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Open response questions

Please write your review here. The author(s) will see this review. Your identity will not be revealed to the authors unless you also include your name (i.e., sign your review) in this box. It is up to you whether to reveal your identity or not, either is fine.

I am the original Reviewer #1. I've read through the revised manuscript and the response letter. I believe that the revision handles most of my original concerns. However, I do still think that the hypotheses could receive greater explanation and justification in the Intro – Table 1 is nice, but ultimately provides just a couple sentences for each hypothesis. I think there should be some additional

explanation/review in the main text. In addition, Table 1 provides an explanation for making that specific prediction, but my original concerns were about both the lack of clarity about the hypotheses (now partially covered by Table 1) and the lack of clarity about the theoretical gap (in my opinion, still missing). Why is it theoretically or practically informative to include SEE and free will in the same experimental paradigm?

The closest I can find is on p. 7: “One important difference is that free will is (mostly) the capacity for action regardless of constraints, both internal and external, and regardless of the outcome, whereas intentionality is a purely internal process, and focused on an association with an outcome. However, these nuanced differences in peoples’ understanding of free will and intention have so far not been comprehensively examined in the literature (Feldman, 2017).” The lack of comprehensive examination on a topic doesn’t mean that we need to examine it. As I said in my previous review, why is it worthwhile to conduct this investigation? Does this identify a bias or inconsistency in how people make these judgments? Can it explain some longstanding philosophical dilemma or persistent real-world phenomena? With that some broader justification like this, I worry that the paper might not receive the impact it could. To be clear: I’m not saying that the questions aren’t worth asking (I think they are), but I think that the authors haven’t done quite enough to justify and sell the questions to a broader audience.

I think the GD does a better job of this to some extent, in particular the discussion of both “bad is stronger than good” and “bad is freer than good” hypotheses. I suggest bringing more of that framing into the introduction.

Author Response

Dec 19, 2024

Reply to Collabra 2nd round decision letter **reviews:** **Free will and side-effect effect**

We would like to thank the editor and the reviewers for their useful suggestions and below we provide a detailed response to each item. We also provide a summary table of changes. Please note that the editor’s and reviewers’ comments are in bold with our reply underneath in normal script.

A track-changes comparison of the previous submission and the revised submission can be found on: <https://draftable.com/compare/heqGIkDqJUar>

A track-changes manuscript is provided with the file:
Collabra-RNR2-FW-SEE-manuscript-v8-G-trackchanges.docx

Reply to Editor: Dr./Prof. Chelsea Helion

I have now read your revised manuscript. I appreciate your careful attention to the concerns the reviewers and I raised. I am happy to provisionally accept your manuscript for submission.

Thank you very much for the positive feedback and the conditional acceptance.

However, the reviewers found a few small things I would like you to address.

Minor revisions:

.1. Reviewer 2 recommends including details about how subjects were recruited (e.g., Mechanical Turk).

In Study 1 we added the clarification that participants were recruited from MTurk using CloudResearch. We also added a note about the filters that we used. These details were already included for Study 2.

.2. Reviewers 2 & 3 noted the helpfulness of Table 1 (I'm noting it as well – it's very clear and significantly enhances readability/comprehension of results!). However, Reviewer 3 (originally Reviewer 1) notes that the Introduction would still benefit from a brief discussion of the theoretical implications that motivated including the side effect effect and free will within the same experimental paradigm.

I think this could possibly be well situated in the “Linking Free Will and SEE: Manipulation of both Agency and Outcome Valence” section, as currently that section explains how you manipulate the variables and that it hasn't been done yet, but not what the theoretical implications would be were you to find what you predict (or its opposite, or a null).

Alternatively, another place this could work well would be where you outline the order of studies under “The Present Investigation” – this section clearly outlines what you did, but not its potential implications for our understanding of either or both phenomena.

Thank you.

We added the following in the subsection “Linking free will and SEE: Free will and intent attributions to side-effects”:

“The side-effect effect paradigm has been used to demonstrate asymmetries in the attribution of intent and blame/praise to seemingly unintended side-effects. The attribution of intent is therefore not only one of the factors associated with the attribution of blame and praise, but, looking at the reverse causal chain, intent is affected by the attribution of blame, so that the need for blame and holding someone accountable leads to stronger attributions of intent (Monroe & Malle, 2019; Malle et al., 2022). Therefore, if even unintended harmful side-effects elicit higher intent, then it is possible that the “bad is freer than good” paradigm identified for free will attributions (Feldman et al., 2016) also extends to lesser chosen or “free” side-effects. That is, if the intentionality side-effect effect extends to free will attributions, then even if the protagonist (e.g., the chairman) only chooses to do something because of focusing on a different unrelated reason (e.g., to increase profits) and that this choice is driven by external pressures (e.g., the board and the shareholders), then with harmful outcomes (e.g., environment is harmed) the protagonist is still attributed as having more free will and the capacity for choice to do otherwise.

Further, examining intent and free will attributions together also helps make clearer the differences between them and their possible links in theories of blame and blame models. For example, in Malle et al. (2014)'s Theory of Blame they provided a “Path Model of Blame” with many different factors, including “intentionality” (“whether the agent brought about the event intentionally”) and “capacity” (“whether the agent could have prevented the event”), yet missing the component of “free will” or “choice” (whether the agent could have chosen whether to prevent the event or not). Choice is loosely related to some of the other factors in the path model, such as “obligation” which serves as an external pressure limiting choice, yet goes far beyond that in capturing internal and

external factors that may have restricted choice (Feldman, 2017). Studying intentionality and free will together by first using two experimental philosophy paradigms that focus on free will and intentionality, and then measuring both free will, intentionality, and blame attributions, can help shed light on 1) the associations between the three and 2) how each factor is affected by manipulations that impact free will and intentionality.”

.3. It’s a bit confusing on p. 15 when it reads “The deterministic and indeterministic universe conditions were presented as follows:” and then has a description of the control condition – I think this may be an error, and you instead want to reference Table 2?

Thank you, we adjusted the description in the method section to make things clearer (underlined was changed/added):

Participants assigned to the deterministic universe and indeterministic universe conditions read a description of the assigned hypothetical universe, then answered comprehension questions and attributions about the described universe to further strengthen the understanding of the described universe. Participants in the universe control condition were not provided with a descriptions of a hypothetical universe. Next, participants were presented with one of the two side-effect effect scenarios. In the deterministic universe and indeterministic universe conditions , the scenarios were described as taking part in the previously described hypothetical universe. The hypothetical universe related descriptions were adjusted from Nichols and Knobe (2007), which contrasted a fully deterministic universe with a universe in which all is deterministic with the exception of humans. In the original study, the two universes were presented together, yet we adjusted the experimental paradigm to split the two descriptions into two different between-subject conditions. The deterministic and indeterministic universe conditions were presented as follows:

Deterministic universe:

Imagine a universe (Universe D) in which everything that happens is completely caused by whatever happened before it. This is true from the very beginning of the universe, so what happened in the beginning of the universe caused what happened next, and so on right up until the present. For example, one-day John decided to have French Fries at lunch. Like everything else, this decision was completely caused by what happened before it. So, if everything in this universe was exactly the same up until John made his decision, then it had to happen that John would decide to have French Fries.

Indeterministic universe:

Imagine a universe (Universe D) in which almost everything that happens is completely caused by whatever happened before it. The one exception is human decision making. For example, one-day John decided to have French Fries at lunch. Since a person's decision in this universe is not completely caused by what happened before it, even if everything in the universe was exactly the same up until John made his decision, it did not have to happen that John would decide to have French Fries. He could have decided to have something different.

[removed previous text]

In Table 2 we adjusted the table note to the following:

Participants in the universe control condition only answered the dependent variables and were not provided with a description of a universe.

.4. On page 40, you note that the correlation ($r = .36$) between free will and blame attributions is strong, but to my understanding, .36 would be indicative of a low-to-moderate correlation – I think you might want to state that it is significant instead?

Thank you. We previously indicated this effect as strong given benchmarks in the social psychology literature such as Richard, Bond Jr., and Stokes-Zoota (2003) indicating correlations of .30 as large effects (see [Jane et al., 2024 Guide to Effect Sizes and Confidence Intervals](#) for more information).

However, this is not an important factor, and more relevant than categorizing the effect strength is the actual effect size already reported there for readers to evaluate, and so we removed the reference to “strong” altogether.

Reply to Reviewer #1

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

Reply to Reviewer #2

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

The authors have done a fine job of responding to my prior concerns. My worry about the discriminant validity of the free will measure remains, but at least the authors acknowledge it in the General Discussion now. Table 1 is exceptionally helpful in summarizing the research and guiding the reader.

Thank you for the positive and supportive opening note and the constructive feedback.

.1. There's only one issue that I noticed. I must have missed this last time, but Study 1 does not contain an explanation of where and how participants were recruited. I'm guessing (based on the reported demographics) that they were recruited via Mechanical Turk, as in Study 2, but this isn't stated explicitly. The authors should add this information.

Thank you.

We revamped our description of the sample in Study 1 to the following:

A total of 427 US American participants were recruited from Amazon Mechanical Turk using CloudResearch (Litman et al., 2017). We employed the following CloudResearch options: Duplicate IP Block, and recruited participants with approval rate of 95% and above and who had more than 100 tasks approved. We first excluded 13 participants who indicated a low English proficiency or self-reported not being serious about filling in the survey. These exclusion criteria were not pre-registered for Study 1, yet we applied it to

be consistent with the pre-registered criteria of Study 2. The exclusion criteria did not have much impact and did not change any of the conclusions of the study (differences in effect size were smaller than 0.1), and we provided the results without exclusions with our code.

We also took this opportunity to add details regarding the participants in Study 2:

We employed the following CloudResearch options: Duplicate IP Block, Block Suspicious Geocode Locations, and Verify Worker Country Location, and recruited participants with approval rate of 95% and above and with 1000-500000 approved tasks.

Reply to Reviewer #3

	Strongly Disagree	Disagree	Neutral	Agree	Strongly Agree
The study/studies in this manuscript have strong construct validity (good measures and/or manipulations of the constructs the authors wish to study). (Choose "Neutral" if this is not an empirical manuscript)				✓	
The study/studies in this manuscript have strong statistical validity (appropriate statistical tests, assumptions are clear and reasonable, no statistical errors, appropriate statistical inferences, etc.). (Choose "Neutral" if this is not an empirical manuscript)					✓
The study/studies in this manuscript have strong internal validity (any causal claims or implications are well-justified, alternative explanations are thoroughly considered, etc.). (Choose "Neutral" if this is not an empirical manuscript, or no causal claims are made or even vaguely implied.)				✓	
The study/studies in this manuscript have strong external validity (authors appropriately constrain their conclusions based on the limits of the generalizability of their findings to other contexts (including from lab to real world), other populations, other stimuli or measures, etc.)				✓	

I am the original Reviewer #1. I've read through the revised manuscript and the response letter. I believe that the revision handles most of my original concerns.

Thank you for the positive and supportive opening note and the constructive feedback.

.1. However, I do still think that the hypotheses could receive greater explanation and justification in the Intro – Table 1 is nice, but ultimately provides just a couple sentences for each hypothesis. I think there should be some additional explanation/review in the main text.

In addition, Table 1 provides an explanation for making that specific prediction, but my original concerns were about both the lack of clarity about the hypotheses (now partially covered by Table 1) and the lack of clarity about the theoretical gap (in my opinion, still missing). Why is it theoretically or practically informative to include SEE and free will in the same experimental paradigm? The closest I can find is on p. 7:

“One important difference is that free will is (mostly) the capacity for action regardless of constraints, both internal and external, and regardless of the outcome, whereas intentionality is a purely internal process, and focused on an association with an outcome. However, these nuanced differences in peoples’ understanding of free will and intention have so far not been comprehensively examined in the literature (Feldman, 2017).”

The lack of comprehensive examination on a topic doesn’t mean that we need to examine it. As I said in my previous review, why is it worthwhile to conduct this investigation? Does this identify a bias or inconsistency in how people make these judgments? Can it explain some longstanding philosophical dilemma or persistent real-world phenomena? With that some broader justification like this, I worry that the paper might not receive the impact it could. To be clear: I’m not saying that the questions aren’t worth asking (I think they are), but I think that the authors haven’t done quite enough to justify and sell the questions to a broader audience.

I think the GD does a better job of this to some extent, in particular the discussion of both “bad is stronger than good” and “bad is freer than good” hypotheses. I suggest bringing more of that framing into the introduction.

Thank you.

We added theoretical explanations in the section “Linking free will and SEE: Free will and intent attributions to side-effects”. See our reply #2 to the editor above.

Editor Final Approval

Jan 4, 2025