

---

# Motivated Political Reasoning: On the Emergence of Belief-Value Constellations

---

**Kai Barron** (WZB Berlin)

**Anna Becker** (Stockholm University)

**Steffen Huck** (University College London, WZB Berlin)

Discussion Paper No. 510

September 16, 2024

# Motivated Political Reasoning: On the Emergence of Belief-Value Constellations\*

Kai Barron<sup>†</sup>, Anna Becker<sup>‡</sup>, Steffen Huck<sup>§</sup>

September 16, 2024

## Abstract

We study the relationship between moral values (“ought” statements) and factual beliefs (“is” statements). We show that thinking about values affects the beliefs people hold. This effect is mediated by prior political leanings, thereby contributing to the polarization of factual beliefs. We document these findings in a pre-registered online experiment with a nationally representative sample of over 1,800 individuals in the US. We also show that participants do not distort their beliefs in response to financial incentives to do so, suggesting that deep values exert a stronger motivational force than financial incentives.

**JEL Codes:** C90, D72, D74, D83, P16

**Keywords:** Motivated Beliefs, Values, Polarization, Experiment, Reasoning.

---

\*The authors would like to thank Ciril Bosch-Rosa, Christoph Drobner, Nathan Hancart, Anna Kerkhof, Nina McMurry, Sebastian Schneider and Alice Solda for helpful comments and discussions. This research was approved by the UCL Research Ethics Committee (17181/001). The study was pre-registered on OSF (<https://osf.io/8jydh/>). Barron gratefully acknowledges financial support from the German Science Foundation via CRC TRR 190 (project number 280092119).

<sup>†</sup>WZB Berlin, Reichpietschufer 50, 10785, Berlin, Germany. Email: kai.barron@wzb.eu

<sup>‡</sup>Stockholm University, 106 91 Stockholm, Sweden. Email: anna.becker@su.se

<sup>§</sup>University College London and WZB Berlin, Reichpietschufer 50, 10785, Berlin, Germany. Email: stef-fen.huck@wzb.eu

# 1 Introduction

Why did Republicans hold different beliefs about the dangers of COVID-19 compared to Democrats, leading them to take fewer precautions and increase their risk of contracting the illness? Both Allcott et al. (2020) and Clinton et al. (2021) document this startling divide in beliefs and behavior between political camps during the pandemic. This example is illustrative of the phenomenon of partisan “bubbles” comprising *disagreement about facts* along the political spectrum. This phenomenon is, however, much wider and can be traced back to, at least, the early 2000s in the United States (see, for example, Gaines et al. 2007, on polarized beliefs about the Iraq war). A common explanation for the emergence and persistence of these patterns is *politically motivated reasoning* (see e.g. Lord, Ross, and Lepper 1979; Taber and Lodge 2006; Kahan 2016). Citizens exhibit a tendency to interrogate arguments and information that conflict with their prior partisan attitudes more vigorously, while uncritically accepting attitudinally congruent arguments.

Politically motivated reasoning has been discussed as a reason for partisan disagreement on key policy issues such as redistribution (Alesina, Stantcheva, and Teso 2018) and immigration (Haaland and Roth 2023; Alesina, Miano, and Stantcheva 2022) and, more generally, as a cause of political polarization (Alesina, Miano, and Stantcheva 2020; Ortoleva and Snowberg 2015). Much of this extant literature is devoted to documenting that beliefs are systematically polarized along political contour lines. However, research on the precise underlying mechanisms that generate these politically polarized beliefs is still scarce. One notable exception to this is Thaler (2024), who presents an innovative experimental study examining how political beliefs may influence the way that individuals draw inference about the veracity of a news source. He shows that individuals tend to trust news more when it aligns with the views of their political party. This reveals one important potential mechanism behind politically polarized beliefs, namely motivated trust in different news sources.

In this study, we examine another important potential mechanism by asking whether the act of *thinking about core values* causally influences the beliefs an individual holds. We conceptualize values as desires about the social world, that is, statements about how the world ought to be (moral positions). Unlike beliefs, which relate to the actual state of the world (how the world is), values cannot be objectively true or false—they can only be endorsed or opposed, with varying degrees of intensity. As previously shown by Alesina, Miano, and Stantcheva (2020) and Enke (2020a), moral values and beliefs are highly correlated. Our study replicates this finding and provides evidence that the salience of values can causally affect factual beliefs. This occurs without the arrival of new information; rather, we examine whether simply encouraging an individual to think about their values can shift their

beliefs. One potential channel for this effect is that making a value salient may distort how an individual retrieves information from their memory to reconstruct their factual beliefs. This distinguishes our study from the belief updating literature which examines distortions in information processing (e.g., Thaler 2024, who uses a clever design where new information is uninformative for a Bayesian, but allows for motivated reasoning by non-Bayesians). Furthermore, we show that individuals distort their beliefs to be more in line with the value system predominant in their political camp. This suggests that political leanings mediate the relationship between values and beliefs that we uncover.<sup>1</sup>

We study the relationship between values and beliefs in the context of six different contentious policy domains: migration, animal welfare, gender equality, abortion, prostitution, and same-sex marriage. In each of these domains, we can test whether individuals adjust their beliefs to conform to the values of their political party. In addition, our experimental design allows us to analyze whether any shift in beliefs translates into a shift in actions. To do this, we give subjects the opportunity to donate money to charities operating in each of the relevant policy domains. Previous studies have shown that politically biased beliefs can have important consequences for individual behavior, for example, in public health (Allcott et al. 2020) or in financial decision-making (Meeuwis et al. 2022).

Our study provides direct evidence on several specific motivated reasoning channels shaping the relationship between values and beliefs and thus also contributes to the broader motivated reasoning literature. In our context, motivated reasoning may operate in two ways. First, since an individual may view her values as an integral part of her identity, motivated reasoning may lead her to shift her factual beliefs so that they support her values. This way, she can maintain a self-image of being someone who holds values that are aligned with facts.<sup>2</sup> Second, when individuals anticipate that holding a certain belief will obligate them to take a costly action (in our context, a donation decision), they may shift their beliefs in the opposite direction to reduce the moral obligation to take that costly action.<sup>3</sup> We examine the role of the second channel by exogenously varying whether individuals anticipate a donation decision. To examine the role of persuasion in political belief formation, we test whether participants shift their beliefs when this can help them to persuade another person to make a certain charitable donation. In doing this, we add to recent findings that demon-

---

1. The importance of conformity with one's preferred political party has previously been shown to be an important factor for politically motivated reasoning. For example, in Druckman, Peterson, and Slothuus (2013) arguments in favor of, or against, a motion are shown to have a stronger effect on partisanship when the arguments are explicitly linked to party stances.

2. For example, an individual who is pro-life might find it easier to believe the results of a study purporting to show that a fetus experiences pain and consciousness than someone who is pro-choice.

3. For example, when an individual anticipates being asked to donate to a charity that supports animal welfare, they might convince themselves that they care a bit less about animal welfare than they usually do.

strate that individuals may shift their own beliefs in order to be more convincing when attempting to persuade another person to believe something (see, e.g., Schwardmann and Van der Weele 2019; Solda et al. 2020; Schwardmann, Tripodi, and Van der Weele 2022).

To explore these questions, we designed and conducted a preregistered online experiment in January 2020 that surveyed a nationally representative sample of 1,863 individuals from the US population.<sup>4</sup> Our experiment employs a between-subject design and comprises four main treatments. These treatments are designed to test hypotheses centered around the following questions: (i) Are there systematic correlations between values and factual beliefs? (ii) Do individuals adjust their factual beliefs when a moral value in the same domain becomes more salient? (iii) Do individuals shift their stated values and factual beliefs to align with their own material self-interest? (iv) Do individuals alter their stated values and beliefs in an attempt to persuade others?

Our first two treatments address questions (i) and (ii) – whether there are correlations between values and beliefs and whether values do indeed causally shape beliefs. First, in treatment `VALUENOTSALIENT`, we elicit subjective beliefs about factual statements from the six domains mentioned above. This treatment serves as a control for the second treatment, `VALUESALIENT`, where we additionally elicit subjects’ agreement with value statements pertaining to the six domains prior to the belief elicitation.

The purpose of this elicitation of values is that it encourages subjects to think about them. This serves to raise the salience of these values when subjects subsequently report their associated beliefs. Of course, subjects may also be passively aware of their values in treatment `VALUENOTSALIENT`, but their direct elicitation in `VALUESALIENT` should serve as a priming device that heightens value salience, bringing values to the forefront of the individual’s mind. The underlying idea behind the `VALUESALIENT` and `VALUENOTSALIENT` treatments is that when confronted with a highly salient value question, individuals may be motivated to shift their factual beliefs to align them with their values.

It is important to note, however, that individuals may already be reporting beliefs that are distorted to align with their values in the `VALUENOTSALIENT` treatment. This implies that what we are able to identify causally is the *additional* shift in beliefs due to exogenously making values more salient. We test this by comparing the distribution of beliefs reported in these two treatments. Since party preferences may mediate the direction of this effect, we examine this comparison when conditioning on party preferences. As outlined in our

---

4. To provide some context, the experiment was, thus, designed and implemented prior to events such as the widespread awareness of COVID-19 (February 2020), the death of George Floyd (May 2020), the claims that the United States presidential election was rigged (November 2020) and the attack on the Capitol (January 2021).

pre-analysis plan, we also report the results for the comparison for the full sample without conditioning.

In all our treatments, there is one final stage after the belief and value elicitation exploring choices that relate to the six domains. Specifically, we give subjects the opportunity to donate money to charities that operate in each of the six domains. This final stage plays a key role for our next two treatments that address questions (iii) and (iv) – whether subjects are willing to distort their values and beliefs in order to convince themselves or to convince others. To test whether self-interest plays a role in shifting beliefs, we introduce our `CONVINCESELF` treatment. This treatment is identical to the `VALUESALIENT` treatment, except that the donation decision is placed on the same screen where we elicit beliefs and values (in other treatments, it comes as a surprise). We can therefore test whether subjects adjust their reasoning in a self-interested way if holding particular belief-value constellations would point towards taking a costly action (i.e., making a donation to a charity whose work is aligned with that particular constellation).

Finally, we test whether subjects adjust their stated beliefs and/or values when they have an incentive to persuade others. In the fourth treatment, `CONVINCEOTHER`, we again ask subjects to state their beliefs and values. However, here, rather than making the donation decision themselves as in the third treatment, they are informed that another participant will have the opportunity to donate after being shown their belief and value responses.

Table 1 provides an overview of the treatments and indicates the number of participants per treatment. We discuss the `BEINGCONVINCED` treatment in more detail at the end of Section 3.1. This treatment is an auxiliary treatment that we implemented to avoid deception. It is not central to our hypotheses or results because participants in this group mostly served as receivers of information from the `CONVINCEOTHER` treatment, which is our primary focus.

Turning to our results, we find support for the existence of aligned belief-value constellations in all the policy domains considered, thereby answering our first research question in the affirmative. Note, however, that while such correlations are evidence of partisanship, it is not possible to understand the mechanism behind this without further evidence. Such bubbles may arise when beliefs shape values, when there are filter bubbles or echo chambers (Flaxman, Goel, and Rao 2016; Enke 2020b), or when values provide a sufficiently strong force for motivated reasoning. It is the comparison of our first two treatments that helps us to explore the latter mechanism.

On aggregate, the distributions of reported subjective beliefs are almost identical across the two treatments. At surface level, this appears to suggest that we do not observe a shift in subjects' beliefs when values are more salient. However, the picture changes dramatically

Table 1: Overview of Experimental Design

	Treatment				
	VALUE-NOTSALIENT	VALUE-SALIENT	CONVINCE-SELF	CONVINCE-OTHER	BEING-CONVINCED
Screen 1	Elicitation of <b>beliefs</b>	Elicitation of <b>values &amp; beliefs</b>	Elicitation of <b>values &amp; beliefs</b> & Option to <b>donate</b>	Elicitation of <b>values &amp; beliefs</b> & Learn about BEINGCONVINCED participant option to <b>donate</b>	Elicitation of <b>values &amp; beliefs</b>
Screen 2	Proceed to next screen as surprise ↓ Option to <b>donate</b>	Proceed to next screen as surprise ↓ Option to <b>donate</b>	End of experiment	End of experiment	Proceed to next screen as surprise ↓ Option to <b>donate</b> (with info. from CONVINCEOTHER participant)
<b>Obs.</b>	<b>375</b>	<b>385</b>	<b>377</b>	<b>363</b>	<b>363</b>

when we control for individuals' political preferences. Specifically, we find that subjects on both the political right and the political left shift their beliefs to align them with the average beliefs held by those in their preferred political camp when values are made more salient. We, thus, do find support for the idea that thinking about values shapes beliefs through motivated reasoning. This suggests that the heightened salience of contentious policy issues in public debates may be a key explanation for increasing polarization in factual beliefs along political attitude division lines.

With respect to our third and fourth research questions, we find that beliefs and values are unaffected by the addition of monetary incentives to persuade oneself or the anticipation

of the opportunity to persuade another person. If anything, this lack of malleability of beliefs and values to other factors appears to suggest that our subjects care about responding honestly to our belief and value questions. This should lend credibility to the internal and external validity of our first two sets of results. These results are also consistent with a growing body of research documenting the limits of motivated reasoning.<sup>5</sup>

*Related Literature:* In studying the role played by values in the (motivated) formation of beliefs, we contribute to a growing literature on motivated cognition and wishful thinking. This literature has considered a wide array of factors that may generate motivated beliefs, including: maintaining a positive image of one's own intelligence, attractiveness, or performance (e.g., Eil and Rao 2011; Möbius et al. 2022; Coutts 2019; Drobner 2022; Huffman, Raymond, and Shvets 2022), judging what is fair or morally appropriate in a self-interested fashion (e.g., Messick and Sentis 1979; Babcock et al. 1995; Konow 2000; Barron, Stüber, and Van Veldhuizen 2019; Amasino, Pace, and Van der Weele 2021), distorting one's own beliefs in order to be more persuasive to others (e.g., Schwardmann and Van der Weele 2019; Solda et al. 2020; Schwardmann, Tripodi, and Van der Weele 2022), and engaging in confirmatory reasoning that reinforces one's prior beliefs (e.g., Nickerson 1998; Rabin and Schrag 1999).

This previous literature typically considers motivated reasoning in relation to a belief that is closely tied to an individual's self-interest, personal characteristics, or pre-existing beliefs about a particular topic. In contrast, here we examine whether deeper values may exert an influence over related factual beliefs. Given the extent to which many contentious political debates are driven by values, along with the substantial heterogeneity in values between and within societies, this strikes us as an important question. Belot and Brisce (2022) provide evidence suggesting that polarization in beliefs can be decreased when individuals are made aware that they share common values such as human rights or behavioral etiquette rules.

---

5. In particular, several of the studies that examine whether belief updating is distorted by monetary incentives associated with different states of the world fail to find any influence of motivated reasoning (see, for example, Gotthard-Real 2017; Coutts 2019; Barron 2021). Furthermore, Thaler (2020) convincingly shows an absence of positivity-motivated reasoning in domains where self-image is not present. A second strand of literature examines scenarios where individuals engage in motivated reasoning to justify or excuse their self-serving behavior in pursuit of monetary gain. This strand of work has revealed mixed results, with Di Tella et al. (2015) and Bicchieri, Dimant, and Sonderegger (2023) documenting evidence that individuals do distort their beliefs to justify self-serving behavior, while the evidence reported in Ging-Jehli, Schneider, and Weber (2020) and Barron, Stüber, and Van Veldhuizen (2019) is more mixed. Nevertheless, while these papers explore the motivated reasoning of individuals seeking monetary gains, the underlying motivation in these studies is linked to beliefs about their character. These motivated beliefs allow individuals to perceive themselves as less of a 'bad person'. Together, these results indicate that motivated reasoning operates in certain domains, with internal psychological factors such as self-image and deep values serving as a source for motivated reasoning, but external factors such as monetary rewards and others' well-being often do not result in motivated reasoning.



Hence, while values might be a divisive factor when they shape belief formation, they also appear to have the potential to reduce gaps between political camps. It is, thus, important to better understand the precise role they play in motivated political reasoning. Our results suggest that when individuals are prompted to think about contentious values, they adjust their factual beliefs to align their beliefs with those of individuals who share their values and political affiliation.

In relation to the *persuasion of others*, in contrast to findings documented in the literature, in our setting, we do not observe evidence that individuals adjust their beliefs to try to be persuasive. One potential explanation for this could be that individuals do not believe that persuasion operates via the transmission of simple statements of beliefs and, rather, requires a richer communication space. For example, individuals may believe that in order to be persuasive, they need to transmit *arguments* (Schwardmann and Van der Weele 2019), *explanations* (Graeber, Roth, and Schesch 2024), or *narratives* (Barron and Fries 2023). Since this richer communication space is ruled out in our design, individuals may not believe that they will be able to persuade others using the limited communication space available and, therefore, not shift their own beliefs.

Our work also contributes empirical evidence to the recent theoretical discussion about how and why partisan individuals increasingly entertain polarized mental models of reality (Leeper and Slothuus 2014; Van Bavel and Pereira 2018; Alesina, Miano, and Stantcheva 2020). One important factor appears to be group identity. As discussed by Sherman and Cohen (2006), individuals tend to select policy-relevant information in a way that allows them to maintain beliefs that are consistent with the position of the group they identify with. The beliefs generated in this way then allow members of identity groups to express behavior that signals group membership and hence strengthens the ties with their own community. Given that liberals and conservatives have been shown to have moral systems based on different psychological foundations (Graham, Haidt, and Nosek 2009), values are likely to be a strong force underlying motivated reasoning. Bonomi, Gennaioli, and Tabellini (2021) formulate a related idea and discuss a theory of identity politics where increasing the salience of a certain policy conflict leads individuals to identify more strongly with their cultural or economic group, and then to distort their beliefs towards the stereotypical belief of the group they identify with. Our results can also be interpreted in the light of this

theoretical framework.<sup>6</sup> Similarly, our work also relates to the large literature investigating how individuals are influenced by social norms, which often results in a desire to conform with the behavior of one’s in-group (see, e.g., Bernheim 1994; Goette, Huffman, and Meier 2006; Bicchieri et al. 2022; Bicchieri, Dimant, and Sonderegger 2023). While this body of work typically focuses on a pull toward conformity in behavior, we focus on the particular mechanism through which beliefs are drawn into tight constellations with members of one’s political in-group.

The remainder of this paper is organized as follows. In Section 2, we describe the experimental setup and results pertaining to questions (i) and (ii) – whether there are correlations between values and beliefs and whether values do indeed causally shape beliefs. Section 3 proceeds with the description of the setup and results for questions (iii) and (iv) – whether subjects are willing to distort their values and beliefs in order to convince themselves or to convince others. Section 4 concludes.

## 2 Existence and Formation of Belief-Value Constellations

Our experimental design consists of four pre-registered<sup>7</sup> treatments that were conducted online using the platform Prolific with a nationally representative sample of 1,863 individuals from the US population.<sup>8</sup> In our main analyses, we essentially follow our pre-registration plan and explicitly flag any minor deviations. In this section we focus on describing and analyzing the first two treatments, VALUE SALIENT and VALUE NOT SALIENT, which allow us to ask: (i) *Do individuals display belief-value constellations? (in the sense of observing a correlation between beliefs and values), and (ii) Do individuals adjust their beliefs to be more coherent with their values?*

---

6. While we focus predominantly on assessing the causal effect of value salience on beliefs, we also contribute to a broader literature that examines the influence of partisanship on information processing. For example, in the domains of energy policy and climate change respectively, Bolsen, Druckman, and Cook (2014) and Druckman and McGrath (2019) examine the role played by partisan differences in information processing due to selectively trusting different *sources* of information. Kahan (2013) explores the role of different thinking styles in generating ideological polarization and Alesina, Stantcheva, and Teso (2018) show that when individuals are provided with pessimistic information about mobility, left-wing individuals become more pessimistic about mobility and increase their demand for redistribution, but right-wing individuals do not. In our paper, individuals are not provided with any new information to process—they must form their beliefs based on the information already stored in their memory. We only vary the presence of a reason for motivated reasoning, such as the salience of a policy conflict.

7. The full pre-registration document can be found at <https://osf.io/8jydh/> and is also reproduced in Appendix D.

8. Table C1 in the appendix shows that our sample is balanced between all the treatments which are described in the following.

## 2.1 The VALUE SALIENT and VALUE NOT SALIENT treatment conditions

The first objective of our VALUE SALIENT treatment is to examine whether we observe a systematic correlation between values and associated factual beliefs. The treatment consists of three parts. First, participants are presented with a sequence of six (randomly ordered) moral value statements and are asked to report the degree to which they agree or disagree with each statement. Table 2 provides an overview of the six different debates together with the moral and factual statements presented to participants. Each of these moral value statements corresponds to a particular contentious topic of debate in public policy, such as gender equality, abortion, or same-sex marriage. The items serve to raise the salience of these value debates such that participants might view later questions through the lens of those debates. It is important to note that we focus on values that lack consensus, which implies that all our topics are inherently political. Hence, raising the salience of a particular value debate may also raise the salience of a political issue. Second, participants are confronted with six factual statements pertaining to the same six domains, and we elicit their beliefs about the veracity of these factual statements.

Both values and beliefs are measured in 5-point Likert scales. For values, responses range from “strongly agree” to “strongly disagree,” and for factual beliefs, from “Very Unlikely” to “Very Likely”.<sup>9</sup> The debates that we consider relate to migration, animal welfare, gender equality, abortion, prostitution, and gay rights. Together, the moral value assessments and factual belief reports allow us to examine whether there is a correlation between individuals’ beliefs and values.

In the last step, we provide participants with the opportunity to make six donation decisions to six separate charities (one of which is randomly implemented). Each of the six charities targets a cause that corresponds to one of the six relevant public policy discussions.<sup>10</sup> For each charity, participants are asked to divide \$3 between the charity and themselves. In a post-experimental survey, we also collected information on the participants’ political attitudes, and we were able to match our data to previously elicited political attitude variables collected by Prolific independently from our experiment.<sup>11</sup> This allows us examine how values and beliefs translate into actions.

---

9. In order to keep the experiment as clear, simple, and easy-to-understand as possible, we opted not to incentivize the elicitation of factual beliefs. This decision draws on the evidence discussed in Haaland, Roth, and Wohlfart (2023) and Stantcheva (2023), indicating that the incentivization of beliefs in this type of study does not necessarily lead to improvements in truth-telling and may distract participants from the key questions of interest.

10. Subjects are provided with information about the aims of the charities and use a slider to indicate how much they would like to donate. Further details about the charities can be found in Table C2 in Appendix D.

11. We provide summary statistics for the main variables from the experiment for all treatments in Table B1 in the Appendix.

Table 2: Overview over Statements and Charities.

	<b>Debate</b>	<b>Moral Statement</b>	<b>Factual Statement</b>	<b>Donation</b>
		<i>“How much do you agree with the following statement?”</i>	<i>“How likely do you think it is that the following statement is true?”</i>	Charity
1	Migration	People should be allowed to migrate freely between countries.	All countries benefit economically from the free movement of labour.	American Immigration Council
2	Animal Welfare	It is wrong to eat animals.	Animals feel less pain than humans.	World Animal Protection
3	Gender Equality	Gender equality should be an objective of policymaking.	Discrimination against women is the primary reason why women earn less than men.	Equality Now
4	Abortion	Abortion should be legal.	Women who have had an abortion experience more psychological distress than women who have had a miscarriage.	Planned Parenthood
5	Prostitution	Prostitution should be illegal.	Human trafficking is facilitated by liberal prostitution laws.	A21
6	Same-sex Marriage	Gay couples should have the same rights as heterosexual couples.	Societies where same-sex marriage is legal are happier than societies where it is illegal.	OutRight

The VALUENOTSALIENT treatment is identical to the VALUE SALIENT treatment, with the exception that the first stage in which participants are presented with moral value statements is skipped. This implies that in this treatment the six public policy debates are not made as salient. The exogenous variation in salience between the two treatments allows us to assess how the shift in salience causally affects factual beliefs. Figures A.1 and A.2 show the instructions as they were presented to participants in both treatments. Importantly, we ask participants about their beliefs concerning factual statements for which there is no clear scientific consensus. We do this because such beliefs are more likely to foster motivated reasoning, as participants do not expect the uncertainty to be resolved. This allows participants to derive anticipatory utility from their beliefs without worrying about learning that they were “wrong.” Drobner (2022) demonstrates that motivated reasoning is more prevalent in

scenarios where individuals do not expect a resolution of uncertainty.

To measure beliefs about factual statements, we ask subjects, “How likely do you think it is that the following statement is true?”; for moral statements, we ask, “How much do you agree with the following statement?” This reflects that for facts, there is, in principle, an ascertainable truth, while values can only be desirable or undesirable to different degrees.

To fix ideas, let’s consider two examples of statement pairs: the first from the migration domain and the second from the animal welfare domain. The statement “All countries benefit from the free movement of labor” pertains to a fact that may be either true or false. While its veracity might be difficult to ascertain, it is, in principle, ascertainable. In contrast, the statement “People should be allowed to migrate freely between countries” expresses a desire. One may or may not agree with the statement, but there is no truth to be ascertained. The same applies to the statements “Animals feel less pain than humans” (belief) and “It is wrong to eat animals” (value).

## 2.2 The Existence of Belief-Value Constellations

The first question we seek to answer is whether there is alignment between the moral values, factual beliefs, and political attitudes that individuals hold. This would indicate the presence of “belief-value constellations”. There are several potential reasons why individuals might have aligned beliefs and values, including: i) holding values that are shaped by beliefs about facts; ii) *avoiding cognitive dissonance* from holding incoherent values and beliefs; or iii) using value and belief statements to justify self-interested actions (i.e., *motivated reasoning*). Our second hypothesis will explore one particular possibility: that beliefs are shaped by values.

Our first set of hypotheses tests whether belief-value constellations are observed systematically in the population.<sup>12</sup>

### **HYPOTHESIS 1: BELIEF-VALUE CONSTELLATIONS**

*There is a correlation between the beliefs, values, and political attitudes that individuals hold. The actions individuals take are aligned with their belief-value constellations.*

*Let  $b_t$  denote the factual beliefs stated by individuals in Treatment  $t \in \{VS, VNS\}$ ,  $v_t$  the moral values stated by individuals,  $d_t$  their donation decisions and  $p_t$  the left-right political stance of individuals.*

---

12. In the interest of facilitating a more coherent exposition of the paper and to enhance readability, we have adjusted the formulation of the hypotheses in comparison to the pre-registration document. We encourage the interested reader to refer to the full pre-registration document in Appendix D for further details.

a) Moral values are positively correlated with beliefs:

$$\text{Corr}(v_{VS}, b_{VS}) \geq 0.$$

b) Moral values are negatively correlated with political attitudes:

$$\text{Corr}(v_{VS}, p_{VS}) \leq 0.$$

c) Donations are positively correlated with beliefs and values:

$$\text{Corr}(d_{VNS}, b_{VNS}) \geq 0,$$

$$\text{Corr}(d_{VS}, b_{VS}) \geq 0, \text{Corr}(d_{VS}, v_{VS}) \geq 0.$$

When reading HYPOTHESIS 1, it is important to take note of the way that the variables are encoded. First, the political stance variables,  $p_t$ , are constructed to be increasing in the degree to which an individual positions herself on the right of the political spectrum. Second, the moral value,  $v_t$ , variables are encoded such that a high value indicates agreement with a value that is typically associated with individuals on the left of the political spectrum. Third, the factual belief variables,  $b_t$ , are defined such that if they are true, they would provide empirical support for moral value positions typically held by individuals on the political left. Finally, the charitable donation variables,  $d_t$ , are constructed such that higher donations are consistent with costly support of a charity aligned with the relevant moral value position.

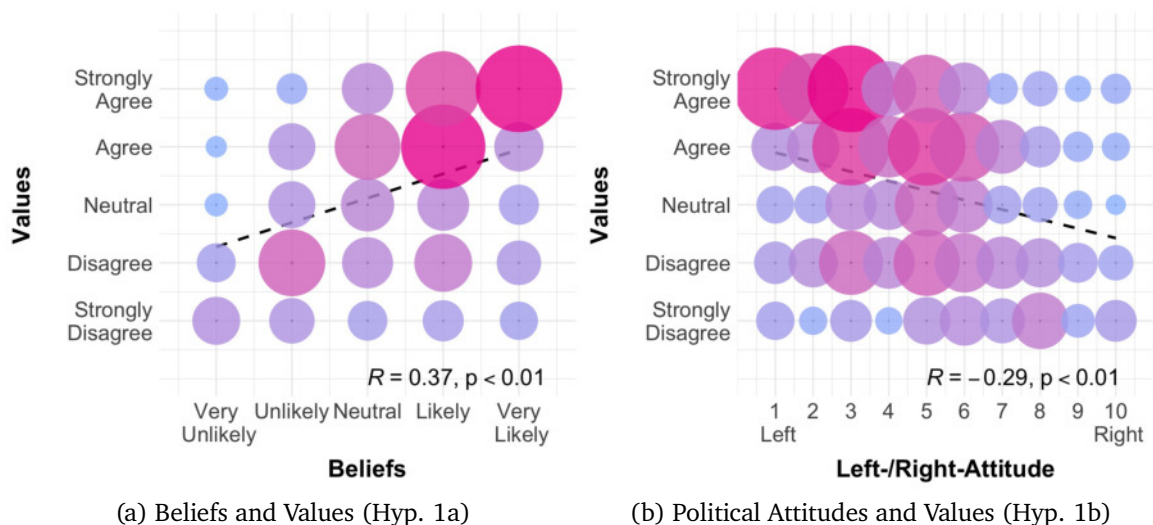
## RESULTS (HYPOTHESIS 1)

Figure 1 summarizes the results pertaining to Hypothesis 1. First, the top left panel reports the correlation between beliefs and values across all policy domains. This shows a strong positive relationship between beliefs and values that is statistically significant at the 1% level. Second, the top right panel shows the results for the correlation between values and political attitudes. In line with the hypothesis, we observe a negative relationship, with left-leaning political attitudes associated with higher agreement with the moral value statements. Third, the three panels in Figure 2 show that donations are positively correlated with beliefs in the VALUENOTSALIENT treatment, and are also positively correlated with beliefs and values in the VALUE SALIENT treatment. All three are statistically significant at the 1% level.

Collectively, these results are in line with the pre-registered set of hypotheses, providing evidence for the presence of belief-value constellations. This suggests that individuals form beliefs, values and political attitudes in a manner that generates strong associations between

the different objects.

Figure 1: Results for Hypothesis 1a and 1b

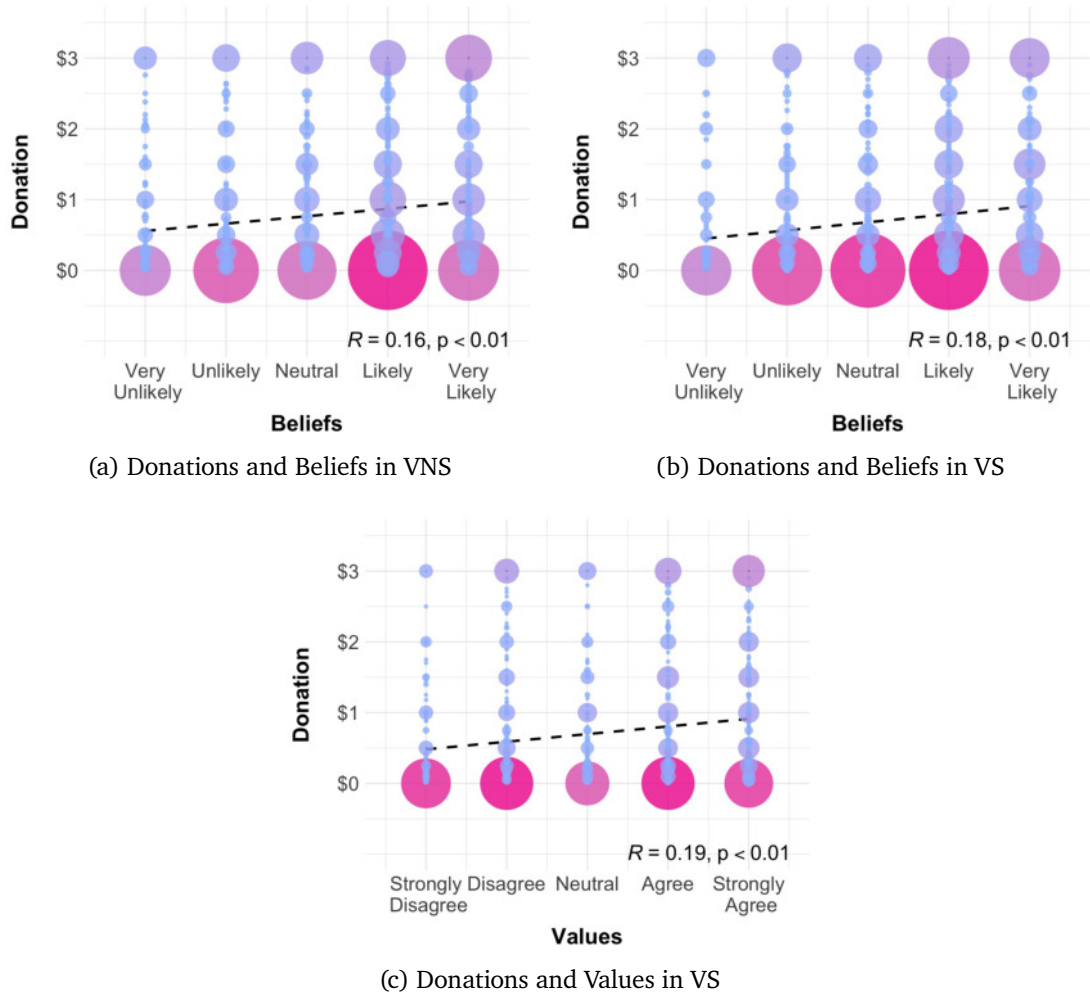


*Note:* Figure 1(a) shows the results on Hypothesis 1a, i.e. the correlation between values and beliefs in the VALUE SALIENT treatment. Figure 1(b) shows the results for Hypothesis 1b, i.e. the correlation between moral values and political attitudes in the VALUE SALIENT treatment. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of values on beliefs in Figure 1(a) and a regression of values on political attitudes in Figure 1(b). The Spearman correlation coefficient,  $R$ , and its  $p$ -value are reported at the bottom right of each graph.

In Appendix B.2, we reproduce these results for each of the six topics separately. Unlike the other five topics, the prostitution topic displays a statistically significant *positive* relationship between political attitudes and values, with individuals who identify with the political right stating stronger agreement that prostitution should be illegal than individuals on the political left.<sup>13</sup> Similarly, figures B.6, B.7 and B.8 in the Appendix illustrate the relationship between donation decisions and beliefs and values in the VALUE NOT SALIENT and VALUE SALIENT treatments for each topic separately. Overall, the relationship between donations and both values and beliefs appears to be weakest for the prostitution-related charity, which received relatively high donation levels across all beliefs and values. Interestingly, for the abortion-related charity, the relationship was very weak in the VALUE NOT SALIENT treatment, but very strong when the value debate was made salient in the VALUE SALIENT treatment.

13. When we designed the experiment, we were aware that the prostitution domain was different to the other domains in the sense that the alignment of aggregate values and political slant was less clear-cut from an ex ante perspective. For this reason, in footnote 35 in Section D.2.1 of our pre-registration document, we noted that our predictions regarding the relationship between values and political slant in the domain of prostitution were more ambiguous.

Figure 2: Results for Hypothesis 1c



*Note:* The three figures show the results pertaining to Hypothesis 1c. Figure 2(a) shows the correlation between beliefs and donations in treatment VALUENOTSALIENT, Figure 2(b) shows the correlation between beliefs and donations in treatment VALUESALIENT, and Figure 2(c) shows the correlation between values and donations in treatment VALUESALIENT. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations (interested readers can also refer to Table B1 in the Appendices for additional descriptive statistics). The dotted line represents the result of a linear regression of donations on beliefs in Figures 2(a) and 2(b), and a regression of donations of values in Figure 2(c). The Spearman correlation coefficient,  $R$ , and its  $p$ -value are given at the bottom right of each graph.

### 2.3 The Formation of Belief-Value Constellations

Our second hypothesis asks whether the formation of factual beliefs is influenced by the salience of a particular contentious moral value debate. *Do individuals adjust their factual beliefs when examining them with a related hotly contested moral issue at the forefront of their mind?* If this is the case, it would speak to a neglected mechanism generating tightly clustered beliefs and values.



To test this, we use a between-treatment comparison of the distribution of beliefs observed in the VALUENOTSALIENT and VALUESALIENT conditions. We can thus assess whether factual beliefs are shifted when we prime individuals to think about these belief statements through the lens of the related value debate.

This is formalized in Hypothesis 2 below, which posits that: (i) the salience of values affects belief formation, and (ii) this mechanism can result in the polarization of factual beliefs. The rationale for this is that if (i) is true then the heterogeneity in moral values between different political groups would also lead to the formation of polarized factual beliefs. This would provide one potential explanation for recently observed trends of increasingly polarized factual beliefs along ideological lines (see, e.g., Gentzkow 2016; Enke 2020a) which has been documented in various domains, such as *climate change* (McCright and Dunlap 2011) and *COVID-19 beliefs* (Allcott et al. 2020).<sup>14</sup>

## **HYPOTHESIS 2: CONSTRUCTION OF CONSISTENT BELIEFS**

*Increasing the salience of a contentious moral value debate leads individuals to report factual beliefs that are more strongly aligned with their moral value position. This results in an increase in polarization of factual beliefs.*

Let  $F_{b_t}$  denote the cumulative distribution function (cdf) of factual beliefs  $b_t$  in treatment  $t \in \{VS, VNS\}$ ,  $F_{v_t}$  the cdf of moral values  $v_t$ , and  $F_{d_t}$  the cdf of donations  $d_t$ . As before,  $p_t$  denotes the left-right political stance of individuals.

a) *Raising the salience of a moral value debate influences factual beliefs.*

*The distribution of factual beliefs differs between the VALUENOTSALIENT and VALUESALIENT treatments:*

$$F_{b_{VNS}} \neq F_{b_{VS}}.$$

b) *A higher degree of polarization in values in a particular policy domain will result in a stronger effect of increased value salience on the dispersion of factual beliefs.*

*Comparing across the six domains indexed by  $m$ , the difference between the variance in beliefs in VALUESALIENT and the variance in beliefs in VALUENOTSALIENT is non-decreasing in the variance in values in VALUESALIENT:*

$$\frac{d[\text{Var}(b_{VS}^m) - \text{Var}(b_{VNS}^m)]}{d[\text{Var}(v_{VS}^m)]} \geq 0.$$

---

14. In his theoretical work, Le Yaouanq (2023) links heterogeneity in political attitudes to partisan disagreement about objective facts through people's idiosyncratic preferences regarding the policy implications of scientific findings. Our work seeks to understand the underlying psychological mechanisms in more detail.

c) *Raising the salience of a moral value debate results in an increase in the polarization of factual beliefs conditional on political attitudes.*

*Conditional on political attitudes, beliefs in VALUE SALIENT are more polarized than beliefs in VALUE NOT SALIENT:*

$$\begin{aligned}
 & E(b_{VS}|p_{VS} < E(p_{VS})) - E(b_{VNS}|p_{VNS} < E(p_{VNS})) \\
 & \geq \\
 & E(b_{VS}|p_{VS} > E(p_{VS})) - E(b_{VNS}|p_{VNS} > E(p_{VNS})).
 \end{aligned}$$

Several features of this set of hypotheses are worth highlighting. First, the rationale for part b) and c) of the hypothesis is that the raised salience of the relevant value will result in a shift towards more extreme factual beliefs as subjects are drawn towards more coherent belief-value constructions. Second, the inequality in part c) states that the difference between the factual beliefs of individuals on the left of the political spectrum and those on the right of the political spectrum will increase in VALUE SALIENT versus VALUE NOT SALIENT (i.e., the salience of values will increase polarization of factual beliefs, conditional on political attitudes).<sup>15</sup>

## RESULTS (HYPOTHESIS 2)

In this section we examine the relationship between values and beliefs more closely by asking whether raising the salience of a particular value leads to a causal shift in an associated factual belief. This comparison represents a fairly conservative test of the existence of a causal relationship between beliefs and values for several reasons. First, our experiment focuses on short-run motivated reasoning and does not consider causal effects of motivated cognition that operate over a longer period of time (e.g., via biased information search or selective memory). Second, our experiment exploits a salience manipulation of values, which represents a fairly weak treatment dosage in relation to an exogenous shift of values. We only encourage participants to *think about* the values they already hold; we do not exogenously shift their values.

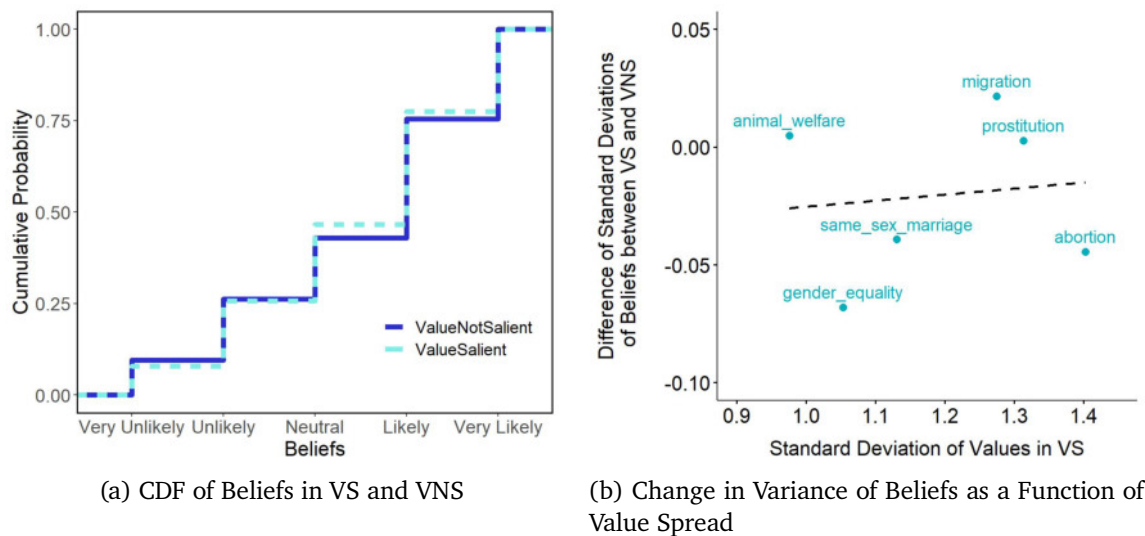
Figure 3 provides a summary of the results associated with Hypothesis 2a and 2b. The left panel displays the cumulative distribution of reported beliefs in the VALUE NOT SALIENT and VALUE SALIENT treatments. We find no statistically significant difference between the two distributions (p-value = 0.09, Kolmogorov-Smirnoff test) and, therefore, do not find support for Hypothesis 2a. Second, the right panel of the figure asks whether there is a heterogeneous effect of increasing the salience of a particular topic. For topics with a high degree

---

15. This can occur due to individuals on the left of the political spectrum increasing their factual beliefs in VALUE SALIENT versus VALUE NOT SALIENT more than those on the right of the political spectrum. Or, it can occur due to individuals on the left increasing their beliefs, while those on the right adjust their beliefs downwards.

of variance in values (i.e. highly polarized issues), we hypothesized that increasing the salience of these values would lead to a larger degree of polarization in the VALUE\_SALIENT beliefs relative to the beliefs in VALUE\_NOT\_SALIENT (Hypothesis 2b). Again, we do not find support for our hypothesis, since the slope coefficient is not statistically different from 0 (coefficient: 0.034; p-value: 0.759).

Figure 3: Results for Hypotheses 2a and 2b



*Note:* Figure 3(a) shows the results for Hypothesis 2a, i.e. the cumulative density function of beliefs in treatments VALUE\_SALIENT and VALUE\_NOT\_SALIENT. Figure 3(b) shows the result for Hypothesis 2b. The y-axis shows the difference of the standard deviations of beliefs between treatments VALUE\_SALIENT and VALUE\_NOT\_SALIENT and the x-axis shows the standard deviation of values in treatment VALUE\_SALIENT. The dotted line depicts the result from a linear regression of the difference of the standard deviations on the standard deviation of values. The slope coefficient is 0.034 (p-value: 0.759).

The results for Hypotheses 2a and 2b suggest that raising the salience of values did not result in a clear shift in the distributions of beliefs across all six issues. Hypothesis 2c posits that even if an increase in value salience does not result in an increase in polarization of the aggregate distribution of beliefs, there may be heterogeneity in the impact of the value salience at the individual level—i.e., the political preferences of an individual could mediate how increasing the salience of their values shifts their beliefs. Essentially, Hypothesis 2c asserts that making a value more salient leads individuals to shift their beliefs even further towards conforming with the average beliefs held by members of their own political party.

To address this question, we compare the belief movement of individuals on the left of the political attitude spectrum with those on the right of the political attitude spectrum. Using a difference-in-difference style empirical approach, we ask whether the gap between the beliefs of those on the left and the right increases in the VALUE\_SALIENT treatment relative

to in the VALUENOTSALIENT treatment.<sup>16</sup> To do this, we estimate the following regression:

$$b_{i,j} = \alpha_0 + \alpha_1 \cdot \tilde{p}_{i,j} + \alpha_2 \cdot ValSal_{i,j} + \alpha_3 \cdot \tilde{p}_{i,j} \times ValSal_{i,j} + \epsilon_{i,j} \quad (1)$$

where  $b_{i,j}$  is the reported belief of individual  $i$  for topic  $j$ ,  $ValSal_{i,j}$  is a binary variable that equals 1 if the individual is in the VALUESALIENT treatment and 0 when the individual is in the VALUENOTSALIENT treatment, and  $\tilde{p}_{i,j}$  is an indicator variable that takes a value of 1 if the individual is on the left of the political spectrum, i.e. reports a political attitude that is lower than the mean political attitude reported in our sample. For this purpose, we asked participants the following question after they completed the experimental part: “In political matters, people talk of “the left” and “the right”. How would you place your views on this scale?”. Respondents could choose a number between 1 and 10, where 1 indicates the extreme left and 10 indicates the extreme right.<sup>17</sup>

The coefficient of interest is  $\alpha_3$ , corresponding to the interaction term. This essentially compares how individuals on the left and right of the political spectrum change their factual beliefs when exposed to an increase in value salience. A positive coefficient denotes a widening of the gap between the factual beliefs of the left and the right. Table 3 reports the results from estimating equation 1 in the first column, with  $VALUESALIENT \times Pol. Attitude$  denoting the interaction term. The estimates show that the coefficient on the interaction term is positive and statistically significant at the one-percent level, providing evidence that we do indeed observe polarization of the beliefs along political attitude division lines when related contentious values are made salient. It is worth noting that this increase in polarization is on top of the pre-existing difference in factual beliefs reported between individuals on the left and right of the political spectrum in the VALUENOTSALIENT treatment. This is shown by the significant coefficient associated with the *Pol. Attitude* variable. It is also worth noting that the size of the widening of the gap in factual beliefs between the left and right due to the salience is nearly as large as the baseline difference in factual beliefs between individuals on the left and the right in VALUENOTSALIENT (i.e. the magnitude of the coefficient associated with the variable  $VALUESALIENT \times Pol. Attitude$  is  $\frac{3}{4}$  the size of the coefficient associated with the variable *Pol. Attitude*).

---

16. In the pre-registration plan, we did not specifically outline this regression, which allows us to test Hypothesis 2c.

17. As a caveat, it should be noted that this question refers to a general political attitude of the respondent and does not elicit their views on social, economic, or other matters separately.

Table 3: Influence of increased salience of values on belief polarization

	(1)	(2)	(3)	(4)	(5)	(6)
VALUESALIENT	-0.124** [0.052]	-0.199** [0.087]	-0.236*** [0.075]	-0.133** [0.052]	-0.222** [0.090]	-0.243*** [0.075]
Pol. Attitude ( $\tilde{p}$ )	0.269*** [0.047]	0.358*** [0.071]	0.304*** [0.058]	0.243*** [0.048]	0.328*** [0.076]	0.284*** [0.061]
<b>VALUESALIENT</b> <b>× Pol. Attitude</b>	<b>0.199***</b> <b>[0.067]</b>	<b>0.294***</b> <b>[0.100]</b>	<b>0.315***</b> <b>[0.089]</b>	<b>0.218***</b> <b>[0.067]</b>	<b>0.316***</b> <b>[0.103]</b>	<b>0.314***</b> <b>[0.089]</b>
Constant	3.325*** [0.038]	3.202*** [0.063]	3.264*** [0.049]	3.428*** [0.070]	3.270*** [0.101]	3.247*** [0.090]
Observations	4560	2550	3006	4548	2544	3000
Incl. Controls	No	No	No	Yes	Yes	Yes
Pol. Attitude ( $\tilde{p}$ )	Left-Right Scale	Party Affiliation	Last Election	Left-Right Scale	Party Affiliation	Last Election
Variable	Left=1	Democrat=1	Clinton=1	Left=1	Democrat=1	Clinton=1

Notes: (i) Each of the regressions uses the observations from the two treatments, VALUENOTSALIENT (375 observations) and VALUESALIENT (385 observations), pooled over all six debates ( $760 \times 6 = 4560$ ). (ii) Smaller sample sizes in columns (2) to (6) result from missing information in the political attitude variables and/or in the control variables. (iii) VALUESALIENT is a dummy variable equal to one if the individual was assigned to treatment VALUESALIENT and hence equal to zero if the individual was assigned to treatment VALUENOTSALIENT. (iv) We use three measures of the political attitudes variable. This is indicated in the last two rows of the table. In column (1), Political Attitude is a dummy equal to one if the individual is below the median on a 1 to 10 scale of political attitudes where 1 is the most left, and 10 is the most right attitude. In column (2), Political Attitude is a dummy equal to one if the individual identifies as a Democrat rather than as a Republican, and in column (3), Political Attitude equals one if the individual indicated that they voted for Clinton in the 2016 elections and zero if they voted for Trump. (v) Columns (4) to (6) show the regressions including controls for age, gender, ethnicity, and education. We present the results for the regression run in column (1) using the subsamples from columns (2) and (3) in Table B2 in the Appendix. (vi) Standard errors clustered at the level of the individual are reported in parentheses, \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

In order to test the robustness of this result, we check whether the results are driven by the specific political attitudes variable that we have chosen to use.<sup>18</sup> To do this, we run two further regressions, where we replicate the estimation in the first column of Table 3, but replace the *Left* indicator variable with a variable that indicates that the individual self-reported being a *Democrat* (Column (2)) and a variable that indicates that the individual voted for Hilary Clinton in 2016 (Column (3)). The results from both of these exercises are

18. In the pre-registration plan, we only mentioned using the political left-right scale as a variable. To test the robustness of the results, we also examine party affiliation and voting decisions in the last election. Additionally, the pre-registration plan does not mention the use of control variables which were also added to test for robustness here.

highly consistent with our main estimation results in Column (1).<sup>19</sup>

These results highlight an important distinction between two forms of polarization, namely (i) polarization of the entire unconditional distribution, which involves movement towards extreme beliefs, and (ii) polarization conditional on a particular characteristic (e.g., political party) that defines groups in the population. The latter form of polarization involves a reshuffling of the belief distribution and may or may not lead to aggregate or unconditional polarization. As seen in Panel (a) in Figure 3, we do not observe unconditional polarization as a result of our experiment, i.e., the distribution over beliefs does not change depending on whether individuals are primed with their values or not. However, this does not reveal whether individuals from different points on the political spectrum have adjusted their beliefs as a result of the value priming. Our regression results in Table 3 point towards important differences between individuals, conditional on their political preferences.

To provide a more detailed explanation, in Figure 4, we disaggregate the beliefs of participants according to their position on the political spectrum.<sup>20</sup> The figure plots the mean beliefs in the two treatments `VALUENOTSALIENT` and `VALUESALIENT` for participants grouped by the political attitude they have indicated on a 10-point scale. Purple stars indicate the means for participants in treatment `VALUENOTSALIENT`, while pink stars indicate the means for individuals in `VALUESALIENT`. We observe that individuals at the extremes of the political spectrum—both on the left and the right—display notable shifts in their mean beliefs when comparing the `VALUESALIENT` condition to the `VALUENOTSALIENT` condition. On the left, the pink stars are above the purple stars; on the right, they are below the purple stars. This asymmetry shows that beliefs move in accordance with political orientation (up on the left; down on the right). In other words, individuals with more extreme political preferences tend to exaggerate their beliefs in a manner that aligns with their group membership when their values are primed.

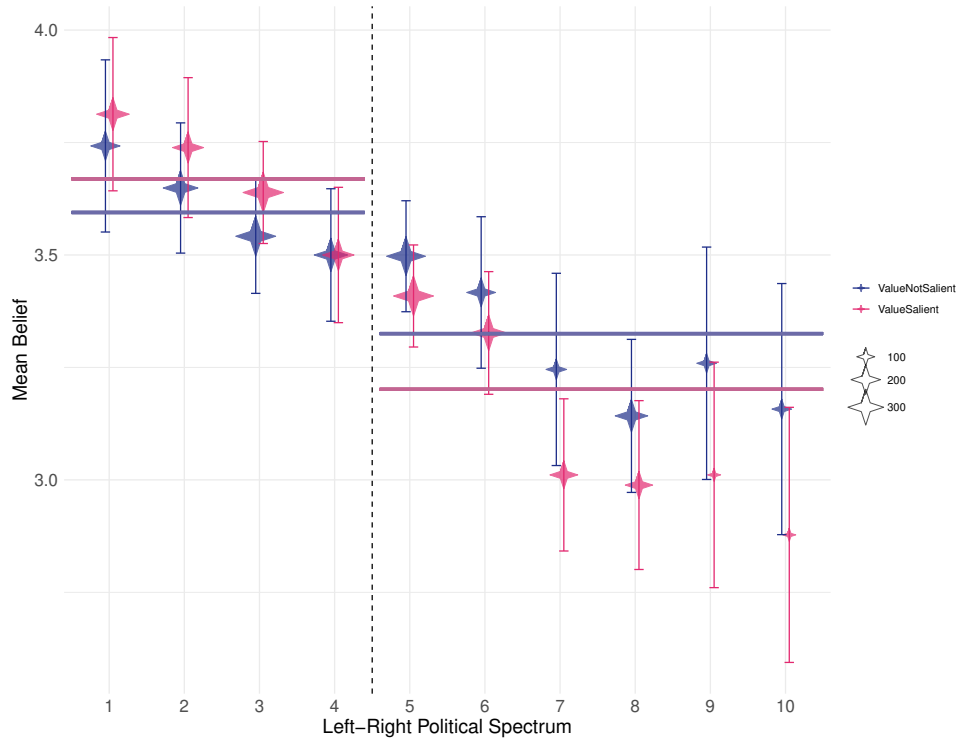
The figure also suggests that the adjustment is larger on the extreme right than on the extreme left, although we need to take into account that we have fewer observations here. In sum, this suggests that while the unconditional distribution of beliefs does not change with the salience of values, the composition of political preferences of people holding different beliefs is affected. Drawing this distinction between *unconditional polarization* and

---

19. Importantly, both of these variables were collected by Prolific completely separately from our experimental data collection. Therefore, these results also serve to alleviate possible concerns regarding our political attitudes variable being influenced by the treatment condition. However, a caveat to this is that the Prolific variables are only available for a subset of the sample. This is the reason for the differing sample sizes across the three regressions.

20. Figure 4 was not mentioned in the pre-registration plan. It was added to enhance the exposition of our results for Hypothesis 2c.

Figure 4: Mean Belief by Political Orientation and Treatment



Note: The stars indicate the mean of beliefs in treatment VALUENOTSALIENT (purple) and VALUESALIENT (pink) along the political attitude of respondents. The associated error bars represent 95% confidence intervals around the means. Responses from the Likert scale regarding beliefs were coded as follows: 1 - “Very Unlikely”, 2 - “Unlikely”, 3 - “Neutral”, 4 - “Likely”, 5 - “Very Likely”. For political attitudes we asked participants the following question after they completed the experimental part: “In political matters, people talk of “the left” and “the right”. How would you place your views on this scale?”. Respondents could choose a number between 1 and 10 where 1 indicates utmost left and 10 indicates utmost right. The vertical dashed line is the mean of political values in the sample comprising observations from Treatment VALUENOTSALIENT and VALUESALIENT. The horizontal bars represent the predicted values from our main regression (see Table 3) for the four groups VALUENOTSALIENT - below the mean (purple line on the left), VALUESALIENT - below the mean (pink line on the left), VALUENOTSALIENT - above the mean (purple line on the right) and VALUESALIENT - above the mean (pink line on the right).

*conditional polarization* is important as it helps us to understand the mechanisms in play. *Unconditional polarization* can be driven by a variety of mechanisms, such as confirmation bias or other individual cognitive heuristics that favor coherent beliefs and values, while *conditional polarization* points towards social conformity with one’s in-group as a driving factor.<sup>21</sup>

### 3 Convincing Yourself and Convincing Others

After having explored how beliefs react to values, we now ask whether money also exerts a distorting influence on beliefs and possibly on values. The third set of our hypotheses below will be divided into two parts, with both parts assessing the malleability of beliefs and values to monetary forces that could pull them in different directions. In Part A (Convincing Yourself), we examine the role of self-serving biases in the context of belief-value constellations by asking whether individuals try to justify selfish behavior by adjusting their beliefs and values to be consistent with taking actions that are in their material self-interest—engaging in a form of motivated reasoning or excuse-driven behavior.<sup>22</sup> Part B (Convincing Others) studies whether introducing the opportunity to try to convince another participant to take an altruistic action can lead to a shift in one’s own beliefs. Specifically, we ask whether attempts to engage in persuasion lead to a shift in beliefs.

#### 3.1 The CONVINCESSELF and CONVINCEOther treatments

We conduct two further treatments. First, the CONVINCESSELF treatment speaks to the conjecture that individuals adjust their beliefs and values in a self-serving way. This treatment is very similar to VALUESALIENT, with only one key difference: in CONVINCESSELF, subjects are aware that they will need to make a charitable donation decision when they form and report their moral value and factual belief assessments. The way that this is implemented in the experiment is that they see the donation decision on the same screen as the one where they report their values and beliefs. This is in contrast to VALUESALIENT, where the charitable donation screen arrives as a surprise after the moral value and factual belief re-

---

21. In Appendix Section B.4.1, we also document the donation behavior in the VALUESALIENT and VALUENOTSALIENT treatment conditions. In summary, we do not observe evidence of a substantial effect on donation decisions, suggesting that the shift in beliefs is not translating into a change in behavior on this dimension. This result contributes to the growing body of work documenting a complex relationship between beliefs and actions.

22. Previous work has shown that people develop self-serving biases in order to excuse their selfishness in charitable giving (see, e.g., Exley (2015) on the role of risk or Exley (2020) on using charity performance metrics as an excuse).



ports have been completed.<sup>23</sup> This difference is important, since subjects' anticipation of the costly charitable donation decision could influence their introspection in forming a personal assessment of the value and factual belief statements. Hypothesis 3a conjectures that individuals bias their (stated) beliefs and values when they take into account the costs of an expected donation decision, with the bias shifting beliefs and values away from those that would justify a higher donation. However, it is also possible that the anticipated donation could have the opposite effect, leading the individual to inflate the importance of her values and beliefs. This self-convincing process may justify a high donation as the correct decision, thereby enhancing her self-image utility from donating.<sup>24</sup>

Second, the CONVINCETHER treatment investigates how trying to persuade others to take an action that is aligned with one's own values could lead an individual to further align their factual beliefs with their political agenda or goals, potentially leading to an exaggeration of these stated beliefs. The underlying assumption here is that subjects get utility from higher donations of others if the donation aligns with their own beliefs and values. To do this, treatment CONVINCETHER mirrors treatment VALUESALIENT with just a single exception: before stating their values and beliefs, subjects are informed that *another* participant will have the option to donate to a related charity after being informed about the moral values and factual beliefs that they (the subject in CONVINCETHER) reported. So, participants might consider the possibility that their own values and beliefs could influence the donation decision of another subject. In order to avoid deception, we implemented these decisions by others in an auxiliary treatment, BEINGCONVINCED. The first part of the BEINGCONVINCED treatment is identical to VALUESALIENT, with subjects reporting their values and beliefs on the six relevant topics. The difference arrives prior to subjects making their donation decisions. At this point, subjects in BEINGCONVINCED are informed about the beliefs and values stated by a randomly chosen participant from CONVINCETHER.

### **3.2 Do individuals self-servingly shift their beliefs and values?**

To examine this question, we compare behavior in CONVINCESSELF, where subjects anticipate their future donation decisions, with behavior in VALUESALIENT, where subjects report their values and beliefs before they are aware of the future donation decisions. This allows us to study the robustness of elicited beliefs and values to the presence of monetary incentives that could distort them. More specifically, we ask whether the presence of the donation decision on the same screen induces subjects to distance themselves from the charity-aligned

---

23. Figures A.1 and A.3 in the Appendix show screenshots of the instructions as they were presented to participants in both treatments.

24. In our pre-registration document we noted this possibility but stated that our prior was that the self-serving bias would dominate.

value position and to adjust their beliefs away from supporting the charity’s goals. This is summarized in the following set of hypotheses.

**HYPOTHESIS 3A: CONVINCING YOURSELF**

As before, let  $F_{b_t}$  denote the cumulative distribution function (cdf) of factual beliefs  $b$  in Treatment  $t$ , and  $F_{v_t}$  the cdf of moral values  $v$ . Donations in Treatment  $t \in \{VS, VNS, CS, CO\}$  are denoted by  $d_t$  and  $p_t$  denotes the left-right political stance of individuals.

- a) *Individuals shift their beliefs and values to justify taking self-serving actions: In CONVINCESelf individuals shift their beliefs and values downwards in comparison to in VALUESalient in order to justify low future donation decisions. Specifically:*
  - i)  $b_{VS}$  first-order stochastically dominates  $b_{CS}$ , i.e.  $F_{b_{VS}} \leq F_{b_{CS}}$ .
  - ii)  $v_{VS}$  first-order stochastically dominates  $v_{CS}$ , i.e.  $F_{v_{VS}} \leq F_{v_{CS}}$ .
- b) *Donations in CONVINCESelf are lower than in VALUESalient:*

$$E(d_2) \geq E(d_3).$$

**RESULTS (HYPOTHESIS 3A)**

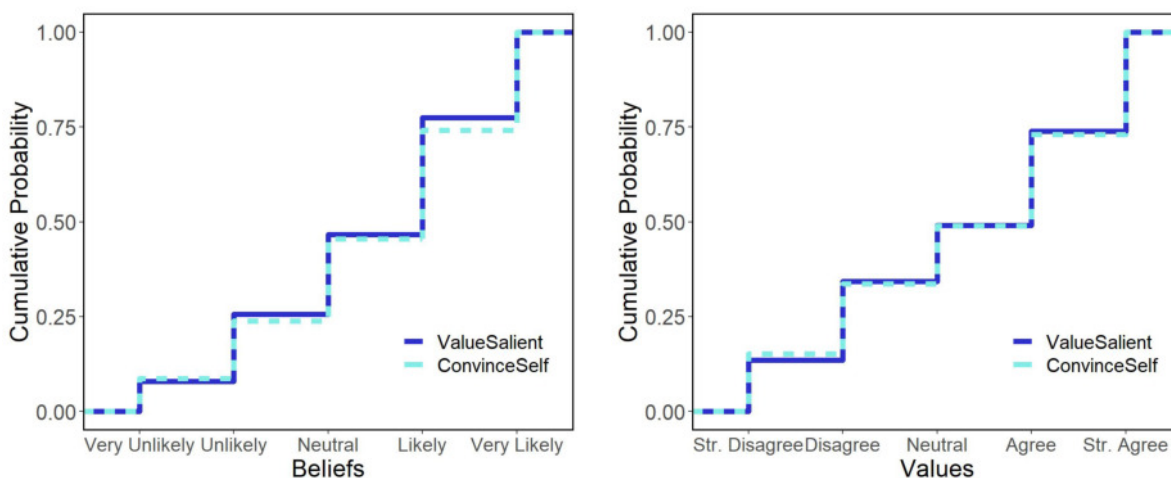
Essentially, we find no evidence in support of Hypothesis 3A. Figure 5 displays the distribution of beliefs (top left panel), values (top right panel) and donations (bottom panel) in the VALUESalient and CONVINCESelf treatments. We observe no significant differences in average behavior between these two treatments, indicating that individuals do not shift their beliefs and values when faced with an imminent donation decision.<sup>25</sup> This immutability of behavior to the anticipated donation decision is in stark contrast to the effect of increased value salience documented above. While subjects are engaging in politically motivated reasoning, they do not engage in economically motivated reasoning. Perhaps one reason for this is that individuals place a higher value on their personal identity, which incorporates their beliefs and values, than they do on a small monetary gain that they would obtain by reducing their donation. A second factor worth noting is that a large fraction of subjects donated less than 1 dollar. Thus, the cognitive dissonance costs of donating a low amount

---

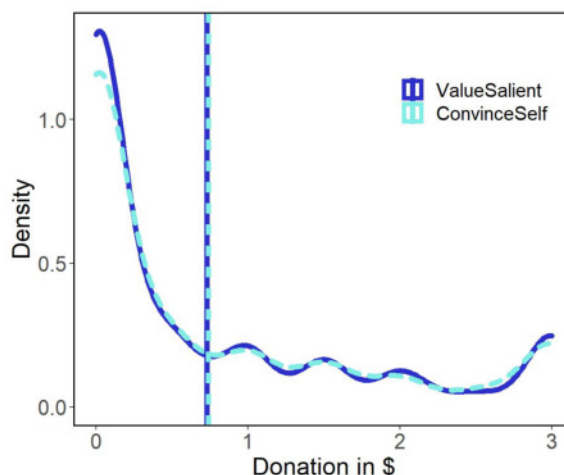
25. We can also separately examine whether donation decisions are affected on either the extensive or intensive margin. Our analysis reveals no significant average differences on either margin. One possible explanation for the failure to observe an average treatment effect is that some individuals shift their beliefs due to monetary self-interest, while others shift their beliefs to enhance the self-image generated from making a donation, with the two effects offsetting each other. Our data does not allow us to convincingly evaluate this possible explanation empirically, but we observe a higher correlation between donations and values and beliefs in the CONVINCESelf treatment in comparison to the VALUESalient treatment. This may be interpreted as evidence in favor of this explanation.

may not be sufficiently high to warrant a shift in beliefs or values to justify it.

Figure 5: Results for Hypothesis 3A



(a) CDF of beliefs in VS and CS (Hypothesis 3A.a.i) (b) CDF of values in VS and CS (Hypothesis 3A.a.ii)



(c) PDF of donations in VS and CS (Hypothesis 3A.b)

*Note:* The three figures show the results on Hypothesis 3. Figure 5(a) shows the cumulative density function of beliefs for treatment VALUE SALIENT (dark line) and CONVINCESelf (light dotted line), Figure 5(b) shows the cumulative density function of values for treatment VALUE SALIENT (dark line) and CONVINCESelf (light dotted line), and Figure 5(c) probability density function of donations for treatment VALUE SALIENT (dark line) and CONVINCESelf (light dotted line). The vertical lines in Figure 5(c) depict the mean of donations in the two treatments.

### 3.3 Do individuals shift their beliefs and values to convince others?

The second part of Hypothesis 3 asks whether individuals report more polarized factual beliefs when they have the opportunity to try to persuade someone else about the importance of certain value positions. It therefore contributes to the body of existing work that examines the idea that we adjust our own beliefs and attitudes (i.e., convince ourselves) in

order to convince others Babcock et al. (1995), Schwardmann and Van der Weele (2019), Solda et al. (2020), and Schwardmann, Tripodi, and Van der Weele (2022). While this previous work predominantly studies scenarios in which an individual is explicitly mandated to convince others about a particular policy position or that they are of high ability, a key difference in our study is that we focus on examining whether individuals try to persuade others to take an action that is aligned with their own values by stating more extreme beliefs. For example, we ask whether an individual might increase their agreement with the statement that “Animals feel less pain than humans.” in order to encourage another person to donate to an animal protection charity.

### **HYPOTHESIS 3B: CONVINCING OTHERS**

*Anticipating the opportunity to persuade another individual about a contentious moral issue shifts one’s own factual beliefs towards the in-group party aligned extreme—i.e., factual beliefs in CONVINCEOTHER are more polarized than factual beliefs in VALUE SALIENT:*

$$\begin{aligned}
 & E(b_{CO}|p_{CO} < E(p_{CO})) - E(b_{VS}|p_{VS} < E(p_{VS})) \\
 & \geq \\
 & E(b_{CO}|p_{CO} > E(p_{CO})) - E(b_{VS}|p_{VS} > E(p_{VS})).
 \end{aligned}$$

Similar to Hypothesis 2c above, the inequality in Hypothesis 3B states that the gap in factual beliefs between individuals on the left and the right of the political spectrum widens when there is an anticipated persuasion opportunity.

### **RESULTS (HYPOTHESIS 3B)**

To examine Hypothesis 3b, Table 4 uses the same empirical specification as above and tests for a divergence of beliefs according to political attitudes between the VALUE SALIENT and the CONVINCEOTHER treatment.<sup>26</sup> Essentially, this asks whether individuals shift their beliefs even further towards conforming with their political in-group when they know that their reports will be viewed by others. The results do not support this hypothesis, with the estimated coefficient on the interaction term close to zero. There are several plausible explanations for this finding, including: (i) that individuals do not wish to persuade others, (ii) that individuals are not prepared to adjust their own beliefs to persuade others, and (iii) that they do not believe that others will be easily persuaded in the context of these contentious debates and the limited communication space available (i.e., not being able to send explanations, narratives or justifications for their persuasive claim). Our data does not

---

26. In the pre-registration plan, we did not specifically outline this regression, which allows us to test Hypothesis 3B. Additionally, the pre-registration plan does not mention the use of control variables which were added to test for robustness here.

allow us to convincingly differentiate between these explanations empirically. Therefore, it is important to keep in mind that multiple possible explanations could account for this result.

Table 4: Incentives to convince others and belief polarization

	(1)	(2)
CONVINCEOTHER	0.089 [0.058]	0.102* [0.056]
Pol. Attitude ( $\tilde{p}$ )	0.468*** [0.048]	0.462*** [0.046]
<b>CONVINCEOTHER <math>\times</math> Pol. Attitude</b>	<b>0.005</b> <b>[0.074]</b>	<b>-0.025</b> <b>[0.072]</b>
Constant	3.201*** [0.036]	3.378*** [0.072]
Observations	4488	4476
Incl. Controls	No	Yes

Standard errors clustered by individual in parentheses.

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

*Note:* Each regression uses the observations from treatments VALUE SALIENT (385 observations) and CONVINCETHOTHER (363 observations), pooled over all six debates. Smaller sample sizes in column (2) result from missing information in the control variables. CONVINCETHOTHER is a dummy variable equal to one if the individual was assigned to treatment CONVINCETHOTHER and hence equal to zero if the individual was assigned to treatment VALUE SALIENT. Political Attitude is a dummy equal to one if the individual is below the median on a 1 to 10 scale of political attitudes where 1 is the most left and 10 is the most right attitude. Column (2) show the regressions including controls for age, gender, ethnicity and education.

## 4 Conclusion

This paper studies whether thinking about moral values can influence beliefs about facts. This is done using a preregistered online experiment that surveyed a nationally representative sample of 1,863 individuals from the US population. We ask two broad sets of questions.

First, we examine whether systematic correlations exist between moral values (“ought” statements) and factual beliefs (“is” statements). We find evidence supporting this relationship, consistent with previous research that highlights a societal shift towards increasingly partisan worldviews (see, e.g., Alesina, Miano, and Stantcheva 2020; Bonomi, Gennaioli, and Tabellini 2021). As we discuss above, there are many mechanisms that might create such a

correlation. For example, beliefs might shape values. Our study examines whether there is a (reverse) causal relationship between values and beliefs, whereby thinking about values exerts an influence on beliefs. We explore this by introducing a treatment that makes a moral value more salient prior to eliciting beliefs. Strikingly, while there appears to be no effect on the aggregate distribution, a closer inspection shows substantial causal effects of thinking about values on beliefs—effects that are mediated by prior political leanings. In other words, we find that individuals in our representative sample are engaged in politically motivated reasoning.<sup>27</sup>

Politically motivated reasoning takes place on both sides of the political spectrum: subjects on both the political right and the political left, shift their beliefs to align them with the average party beliefs when values are made salient. This finding contrasts with the popular belief that the flirtation with “alternative facts” is a phenomenon exclusive to populist right-wing movements.

The influence of *thinking about values* may be generated by different potential psychological mechanisms underlying motivated reasoning. One possibility is that thinking about a particular value position strengthens the individual’s desire for this position to be justified by facts; consequently they shift their beliefs (“pure motivated reasoning”). Another possibility is that thinking about a particular value position cues recall of information in the individual’s memory database that is supportive of the value position (“motivated memory”). This relationship between memory and biased belief has been explored in recent work on associative memory (Enke, Schwerter, and Zimmermann 2024) and motivated memory (Zimmermann 2020; Amelio and Zimmermann 2023). Our experiment is not able to cleanly distinguish between these channels through which the effect of thinking about values may operate. We leave this interesting question for future work.

Additionally, we examine whether there is evidence for economically motivated reasoning whereby individuals bias their beliefs and/or values due to the presence of monetary incentives to do so. This is not the case. We believe that this result enhances the credibility of our main findings. Since beliefs and values do not react to (small) monetary incentives, it appears that individuals care about them to the extent that they do not distort them through

---

27. The behavior observed in our study is consistent with the findings of Bordalo, Tabellini, and Yang (2021), who study the effect of issue salience on beliefs about others’ political attitudes. The authors show that when the salience of a particular policy conflict is raised, this increases the perception of the partisan gap in attitudes. Combined with an identity-induced desire to conform to the stereotypical beliefs of one’s identity group (as in Bonomi, Gennaioli, and Tabellini 2021), this perceived increase in the partisan gap could contribute to the shift in beliefs that we observe. One caveat is worth keeping in mind when interpreting our results: We focus on a particular kind of belief, namely those without a scientific consensus. For facts where such a scientific consensus does exist, the findings of Drobner (2022) suggest that we might expect less motivated reasoning since individuals will anticipate uncertainty resolution.

economically motivated reasoning.

When interpreting our results, several considerations should be kept in mind. In footnote 9, we explained our rationale for not incentivizing the belief elicitation of factual beliefs. Although evidence from Danz, Vesterlund, and Wilson (2022), Haaland, Roth, and Wohlfart (2020), and Stantcheva (2023) suggests that this decision is unlikely to cause systematic distortions in reported beliefs, it is still possible that individuals report beliefs that differ from those they hold in mind. Additionally, it is possible that experimenter demand effects might influence the beliefs reported by participants in our experiment. While we consider it unlikely that demand effects systematically account for the pattern of beliefs reported across treatments (e.g., we observe no effect of the CONVINCETHE treatment manipulation where demand effects might be expected to be strongest), we cannot completely rule out this possibility. Therefore, these potential influences should be kept in mind when interpreting the results.

Taken together, our results point towards a deep (cognitive) link between values and beliefs. This tight relationship between values and beliefs is consistent with the conceptual idea of a “polarized reality”, where individuals perceive reality through the lens of their economic or social identity (Alesina, Miano, and Stantcheva 2020) and then adjust their beliefs to conform to the stereotypical belief of the salient identity group (Bonomi, Gennaioli, and Tabellini 2021). More broadly, this recent line of research showing how identity shapes beliefs through the desire for group conformity builds on a longer history of research examining how identity can generate a desire for conformity in *actions* (Akerlof and Kranton 2000, 2005; Shayo 2020). With the polarization of social discourse (particularly online) seemingly increasing in society, this body of work points towards identity-induced belief conformity as an important avenue for further research.

Finally, our results also suggest that individuals may engage more readily in politically motivated reasoning than economically motivated reasoning. One caveat to this assertion is that the economic incentives in our experiment are limited. Nevertheless, this points towards a potentially important division of the space of possible motivated reasoning domains. This, too, is an interesting avenue for further research.

## References

- Akerlof, George A., and Rachel E. Kranton. 2000. "Economics and Identity." *Quarterly Journal of Economics* 115 (3): 715–753.
- . 2005. "Identity and the Economics of Organizations." *Journal of Economic Perspectives* 19 (1): 9–32.
- Alesina, Alberto, Armando Miano, and Stefanie Stantcheva. 2020. "The Polarization of Reality." *AEA Papers and Proceedings* 110:324–28.
- . 2022. "Immigration and Redistribution." *The Review of Economic Studies* 90, no. 1 (March): 1–39.
- Alesina, Alberto, Stefanie Stantcheva, and Edoardo Teso. 2018. "Intergenerational Mobility and Preferences for Redistribution." *American Economic Review* 108 (2): 521–54.
- Allcott, Hunt, Levi Boxell, Jacob Conway, Matthew Gentzkow, Michael Thaler, and David Yang. 2020. "Polarization and Public Health: Partisan Differences in Social Distancing During the Coronavirus Pandemic." *Journal of Public Economics* 191:104254.
- Amasino, Dianna, Davide Pace, and Joël van der Weele. 2021. "Fair Shares and Selective Attention." *Working Paper*.
- Amelio, Andrea, and Florian Zimmermann. 2023. "Motivated Memory in Economics—A Review." *Games* 14 (1): 15.
- Babcock, Linda, George Loewenstein, Samuel Issacharoff, and Colin Camerer. 1995. "Biased judgments of fairness in bargaining." *American Economic Review* 85 (5): 1337–1343.
- Barron, Kai. 2021. "Belief Updating: Does the 'good-news, bad-news' Asymmetry Extend to Purely Financial Domains?" *Experimental Economics* 24 (1): 31–58.
- Barron, Kai, and Tilman Fries. 2023. "Narrative Persuasion." *CESifo Working Paper No. 10206*.
- Barron, Kai, Robert Stüber, and Roel van Veldhuizen. 2019. "Motivated Motive Selection in the Lying-dictator Game." *WZB Discussion Paper*.
- Belot, Michele, and Guglielmo Briscese. 2022. "Bridging America's Divide on Abortion, Guns and Immigration: An Experimental Study." *CEPR Discussion Paper 17444*.
- Bernheim, Douglas. 1994. "A Theory of Conformity." *Journal of Political Economy* 102 (5): 841–877.



- Bicchieri, Cristina, Eugen Dimant, Simon Gächter, and Daniele Nosenzo. 2022. "Social Proximity and the Erosion of Norm Compliance." *Games and Economic Behavior* 132:59–72.
- Bicchieri, Cristina, Eugen Dimant, and Silvia Sonderegger. 2023. "It's Not a Lie if You Believe the Norm Does not Apply: Conditional Norm-following and Belief Distortion." *Games and Economic Behavior* 138:321–354.
- Bolsen, Toby, James N. Druckman, and Fay Lomax Cook. 2014. "The Influence of Partisan Motivated Reasoning on Public Opinion." *Political Behavior* 36 (2): 235–262.
- Bonomi, Giampaolo, Nicola Gennaioli, and Guido Tabellini. 2021. "Identity, Beliefs, and Political Conflict." *Quarterly Journal of Economics* 136 (4): 2371–2411.
- Bordalo, Pedro, Marco Tabellini, and David Yang. 2021. "Issue Salience and Political Stereotypes." *NBER Working Paper*.
- Clinton, Joshua, Jon Cohen, John Lapinski, and Marc Trussler. 2021. "Partisan Pandemic: How Partisanship and Public Health Concerns Affect Individuals' Social Mobility during COVID-19." *Science Advances* 7 (2): eabd7204.
- Costa-Gomes, Miguel A., Steffen Huck, and Georg Weizsäcker. 2014. "Beliefs and Actions in the Trust Game: Creating Instrumental Variables to Estimate the Causal Effect." *Games and Economic Behavior* 88:298–309.
- Costa-Gomes, Miguel A., and Georg Weizsäcker. 2008. "Stated Beliefs and Play in Normal-form Games." *Review of Economic Studies* 75 (3): 729–762.
- Coutts, Alexander. 2019. "Good news and bad news are still news: Experimental evidence on belief updating." *Experimental Economics* 22 (2): 369–395.
- Danz, David, Lise Vesterlund, and Alistair J. Wilson. 2022. "Belief Elicitation and Behavioral Incentive Compatibility." *American Economic Review* 112 (9): 2851–2883.
- Di Tella, Rafael, Ricardo Perez-Truglia, Andres Babino, and Mariano Sigman. 2015. "Conveniently upset: Avoiding altruism by distorting beliefs about others' altruism." *American Economic Review* 105 (11): 3416–42.
- Drobner, Christoph. 2022. "Motivated Beliefs and Anticipation of Uncertainty Resolution." *American Economic Review: Insights* 4 (1): 89–105.
- Druckman, James N., and Mary C. McGrath. 2019. "The Evidence for Motivated Reasoning in Climate Change Preference Formation." *Nature Climate Change* 9 (2): 111–119.

- Druckman, James N., Erik Peterson, and Rune Slothuus. 2013. "How Elite Partisan Polarization Affects Public Opinion Formation." *American Political Science Review* 107 (1): 57–79.
- Eil, David, and Justin M Rao. 2011. "The good news-bad news effect: asymmetric processing of objective information about yourself." *American Economic Journal: Microeconomics* 3 (2): 114–38.
- Enke, Benjamin. 2020a. "Moral Values and Voting." *Journal of Political Economy* 128 (10): 3679–3729.
- . 2020b. "What You See Is All There Is." *Quarterly Journal of Economics* 135 (3): 1363–1398.
- Enke, Benjamin, Frederik Schwerter, and Florian Zimmermann. 2024. "Associative Memory, Beliefs and Market Interactions." *Journal of Financial Economics* 157:103853.
- Exley, Christine L. 2015. "Excusing Selfishness in Charitable Giving: The Role of Risk." *Review of Economic Studies* 83, no. 2 (October): 587–628.
- . 2020. "Using Charity Performance Metrics as an Excuse Not to Give." *Management Science* 66 (2): 553–563.
- Flaxman, Seth, Sharad Goel, and Justin M. Rao. 2016. "Filter Bubbles, Echo Chambers, and Online News Consumption." *Public Opinion Quarterly* 80, no. S1 (March): 298–320.
- Gaines, Brian J., James H. Kuklinski, Paul J. Quirk, Buddy Peyton, and Jay Verkuilen. 2007. "Same Facts, Different Interpretations: Partisan Motivation and Opinion on Iraq." *Journal of Politics* 69 (4): 957–974.
- Gentzkow, Matthew. 2016. *Polarization in 2016*. Technical report. Toulouse Network for Information Technology Working Paper.
- Ging-Jehli, Nadja R., Florian H. Schneider, and Roberto A. Weber. 2020. "On Self-serving Strategic Beliefs." *Games and Economic Behavior* 122:341–353.
- Goette, Lorenz, David Huffman, and Stephan Meier. 2006. "The Impact of Group Membership on Cooperation and Norm Enforcement: Evidence Using Random Assignment to Real Social Groups." *American Economic Review* 96 (2): 212–216.
- Gotthard-Real, Alexander. 2017. "Desirability and Information Processing: An Experimental Study." *Economics Letters* 152:96–99.

- Graeber, Thomas, Christopher Roth, and Constantin Schesch. 2024. "Explanations." *ECONtribute Discussion Paper*.
- Graham, Jesse, Jonathan Haidt, and Brian A. Nosek. 2009. "Liberals and Conservatives Rely on Different Sets of Moral Foundations." *Journal of Personality and Social Psychology* 96 (5): 1029–1046.
- Haaland, Ingar, and Christopher Roth. 2023. "Beliefs about Racial Discrimination and Support for Pro-black Policies." *Review of Economics and Statistics* 105 (1): 40–53.
- Haaland, Ingar, Christopher Roth, and Johannes Wohlfart. 2020. "Designing Information Provision Experiments." *CEBI Working Paper Series, Working Paper 20/20*.
- . 2023. "Designing Information Provision Experiments." *Journal of Economic Literature* 61 (1): 3–40.
- Huffman, David, Collin Raymond, and Julia Shvets. 2022. "Persistent Overconfidence and Biased Memory: Evidence from Managers." *American Economic Review* 112, no. 10 (October): 3141–75.
- Ivanov, Asen, Dan Levin, and Muriel Niederle. 2010. "Can Relaxation of Beliefs Rationalize the Winner's Curse?: An Experimental Study." *Econometrica* 78 (4): 1435–1452.
- Kahan, Dan M. 2013. "Ideology, Motivated Reasoning, and Cognitive Reflection: An Experimental Study." *Judgment and Decision Making* 8:407–24.
- . 2016. "The Politically Motivated Reasoning Paradigm, Part 1: What Politically Motivated Reasoning Is and How to Measure It." In *Emerging Trends in the Social and Behavioral Sciences*, 1–16. John Wiley & Sons, Ltd.
- Konow, James. 2000. "Fair Shares: Accountability and Cognitive Dissonance in Allocation Decisions." *American Economic Review* 90 (4): 1072–1091.
- Le Yaouanq, Yves. 2023. "A model of voting with motivated beliefs." *Journal of Economic Behavior & Organization* 213:394–408.
- Leeper, Thomas J., and Rune Slothuus. 2014. "Political Parties, Motivated Reasoning, and Public Opinion Formation." *Political Psychology* 35:129–156.
- Lord, Charles G., Lee Ross, and Mark R. Lepper. 1979. "Biased Assimilation and Attitude Polarization: The Effects of Prior Theories on Subsequently Considered Evidence." *Journal of Personality and Social Psychology* 37 (11): 2098–2109.

- McCright, Aaron M., and Riley E. Dunlap. 2011. "The Politicization of Climate Change and Polarization in the American Public's Views of Global Warming, 2001–2010." *Sociological Quarterly* 52 (2): 155–194.
- Meeuwis, Maarten, Jonathan A. Parker, Antoinette Schoar, and Duncan Simester. 2022. "Belief Disagreement and Portfolio Choice." *The Journal of Finance* 77 (6): 3191–3247.
- Messick, David M., and Keith P. Sentis. 1979. "Fairness and Preference." *Journal of Experimental Social Psychology* 15 (4): 418–434.
- Möbius, Markus M., Muriel Niederle, Paul Niehaus, and Tanya S. Rosenblat. 2022. "Managing Self-Confidence: Theory and Experimental Evidence." *Management Science* 68 (11): 7793–7817.
- Nickerson, Raymond S. 1998. "Confirmation Bias: A Ubiquitous Phenomenon in Many Guises." *Review of General Psychology* 2 (2): 175–220.
- Ortoleva, Pietro, and Erik Snowberg. 2015. "Overconfidence in Political Behavior." *American Economic Review* 105 (2): 504–35.
- Rabin, Matthew, and Joel L. Schrag. 1999. "First Impressions Matter: A Model of Confirmatory Bias." *Quarterly Journal of Economics* 114 (1): 37–82.
- Schwardmann, Peter, Egon Tripodi, and Joël J. van der Weele. 2022. "Self-Persuasion: Evidence from Field Experiments at International Debating Competitions." *American Economic Review* 112 (4): 1118–46.
- Schwardmann, Peter, and Joël van der Weele. 2019. "Deception and Self-deception." *Nature Human Behaviour* 3 (10): 1055–1061.
- Shayo, Moses. 2020. "Social Identity and Economic Policy." *Annual Review of Economics* 12:355–389.
- Sherman, David K., and Geoffrey L. Cohen. 2006. "The Psychology of Self-defense: Self-affirmation Theory." *Advances in Experimental Social Psychology*, 183–242.
- Solda, Alice, Changxia Ke, Lionel Page, and William von Hippel. 2020. "Strategically Delusional." *Experimental Economics* 23 (3): 604–631.
- Stantcheva, Stefanie. 2023. "How to Run Surveys: A Guide to Creating Your Own Identifying Variation and Revealing the Invisible." *Annual Review of Economics* 15:205–234.

- Taber, Charles S., and Milton Lodge. 2006. "Motivated Skepticism in the Evaluation of Political Beliefs." *American Journal of Political Science* 50 (3): 755–769.
- Thaler, Michael. 2020. "Do People Engage in Motivated Reasoning to Think the World Is a Good Place for Others?" *arXiv preprint arXiv:2012.01548*.
- . 2024. "The Fake News Effect: Experimentally Identifying Motivated Reasoning Using Trust in News." *American Economic Journal: Microeconomics* 16 (2): 1–38.
- van Bavel, Jay J., and Andrea Pereira. 2018. "The Partisan Brain: An Identity-based Model of Political Belief." *Trends in Cognitive Sciences* 22 (3): 213–224.
- Zimmermann, Florian. 2020. "The Dynamics of Motivated Beliefs." *American Economic Review* 110 (2): 337–363.

# A Experimental Design

Figure A.1: Starting Screens

Welcome!

This survey is conducted by researchers from University College London (UCL) and the Berlin Social Science Center (WZB).

Please read the following carefully before you decide whether you consent to participate.

You will be asked questions on a broad range of different topics. Among those will be topics you might perceive as sensitive (for example on your political opinions). Please note that all data from this experiment will be completely anonymised. It will not be possible to link the information you provide to any information that would allow to identify you. The data will only be used for the purpose of this study. If you would like to withdraw your consent, please get in touch with Prolific.

Check the following box if you consent to participate in this study:

Please pay close attention to all questions. We are interested in your opinion/best guess only so please do not try to find the answers to the questions we ask you online.

Please click below if you are ready to begin.

[Next](#)

(a) VALUENOTSALIENT and VALUESALIENT

Welcome!

This survey is conducted by researchers from University College London (UCL) and the Berlin Social Science Center (WZB).

Please read the following carefully before you decide whether you consent to participate.

You will be asked questions on a broad range of different topics. Among those will be topics you might perceive as sensitive (for example on your political opinions). Please note that all data from this experiment will be completely anonymised. It will not be possible to link the information you provide to any information that would allow to identify you. The data will only be used for the purpose of this study. If you would like to withdraw your consent, please get in touch with Prolific.

Check the following box if you consent to participate in this study:

Please pay close attention to all questions. We are interested in your opinion/best guess only so please do not try to find the answers to the questions we ask you online.

In the first part of this survey you will have the **option** to donate to charity six times. One of your decisions will be then chosen at random. For this decision, we will donate the amount you have chosen to the charity and pay the rest as a bonus to you at the end of the survey. Note that the rules of our institute do not permit deception of participants, so all promised payments and donations will actually be made after the experiment. As this study is anonymous the donation cannot be linked to you.

Please click below if you are ready to begin.

[Next](#)

(b) CONVINCESSELF

Welcome!

This survey is conducted by researchers from University College London (UCL) and the Berlin Social Science Center (WZB).

Please read the following carefully before you decide whether you consent to participate.

You will be asked questions on a broad range of different topics. Among those will be topics you might perceive as sensitive (for example on your political opinions). Please note that all data from this experiment will be completely anonymised. It will not be possible to link the information you provide to any information that would allow to identify you. The data will only be used for the purpose of this study. If you would like to withdraw your consent, please get in touch with Prolific.

Check the following box if you consent to participate in this study:

Please pay close attention to all questions. We are interested in your opinion/best guess only so please do not try to find the answers to the questions we ask you online.

**After you have answered all questions another participant of the experiment will be chosen randomly and have the option to make a donation to a charity as you will see it on your screen.** One of their decisions will be then chosen at random. For this decision, we will donate the amount they have chosen to the charity and pay the rest as a bonus to them at the end of the survey. Note that the rules of our institute do not permit deception of participants, so all promised payments and donations will actually be made after the experiment. As this study is anonymous the donation cannot be linked to you.

Please click below if you are ready to begin.

[Next](#)

(c) CONVINCETHER

Note: Panels (a) to (c) show the starting screens for participants in the respective treatments.

Figure A.2: Decision Screens Treatment VALUE SALIENT and VALUE NOT SALIENT

(a) First Decision in VALUE NOT SALIENT

(b) First Decision in VALUE SALIENT

(c) Instructions in VALUE NOT SALIENT and VALUE SALIENT

(d) Second Decision in VALUE NOT SALIENT

(e) Second Decision in VALUE SALIENT

Note: Participants in treatments VALUE NOT SALIENT and VALUE SALIENT were first presented with the screens in Panel (a) and Panel (b), respectively. Afterwards, they were presented with the information in Panel (c), before they would take their decisions as instructed in Panel (d) and (e), respectively. **Please note that the debates participants were presented with were the same in the first and in the second decision in the experiment.** Participants could proceed to the next screen after the first and second decision after waiting for 15 seconds. The starting value of the slider was randomly chosen by the computer.

Figure A.3: Decision Screens Treatment CONVINCESelf and CONVINCeOther

How much do you agree or disagree with the following statement?

**Abortion should be legal.**

Strongly Disagree  Disagree  Neutral  Agree  Strongly Agree

How likely do you think it is that the following statement is true?

**Women who had an abortion experience more psychological distress than women who had a miscarriage.**

Very Unlikely  Unlikely  Neutral  Likely  Very Likely

You now have the opportunity to donate to **Planned Parenthood**. The mission of Planned Parenthood is to provide comprehensive reproductive and complementary health care services and to advocate public policies which ensure access to such services. They provide information and support to women considering to end a pregnancy in a clinic or using an abortion pill.

You have an amount of \$3 available. Your bonus will be determined as \$3 minus what you donate.

Use the slider to indicate how much would you like to donate.

Amount that will be donated to Planned Parenthood: **\$1.4**

Please read all statements carefully. You can continue to the next page at any time after 9 seconds.

(a) CONVINCESelf

How much do you agree or disagree with the following statement?

**Gay couples should have the same rights as heterosexual couples.**

Strongly Disagree  Disagree  Neutral  Agree  Strongly Agree

How likely do you think it is that the following statement is true?

**Societies where same-sex marriage is legal are happier than societies where it is illegal.**

Very Unlikely  Unlikely  Neutral  Likely  Very Likely

Another participant will face the following decision after being informed about your answer on the question above.

You now have the opportunity to donate to **Outright**. They envision a world where LGBTQ people everywhere enjoy full human rights and fundamental freedoms. They seek to fill research gaps, provide trainings to community members and allies to develop their expertise, and convene key stakeholders to exchange information on best practices related to ending violence based on sexual orientation.

You have an amount of \$3 available. Your bonus will be determined as \$3 minus what you donate.

Use the slider to indicate how much would you like to donate.

Amount that will be donated to Outright: **\$3**

Please read all statements carefully. You can continue to the next page at any time after 6 seconds.

(b) CONVINCeOther

Note: Participants in treatments CONVINCESelf and CONVINCeOther were presented with the screens in Panel (a) and Panel (b), respectively. Participants could proceed to the next screen after waiting for 15 seconds. The starting value of the slider was randomly chosen by the computer. In CONVINCeOther, they could not interact with the slider and were presented with the screen that the other participant would be shown.



## B Supplementary Results

### B.1 Descriptive Statistics for Main Variables

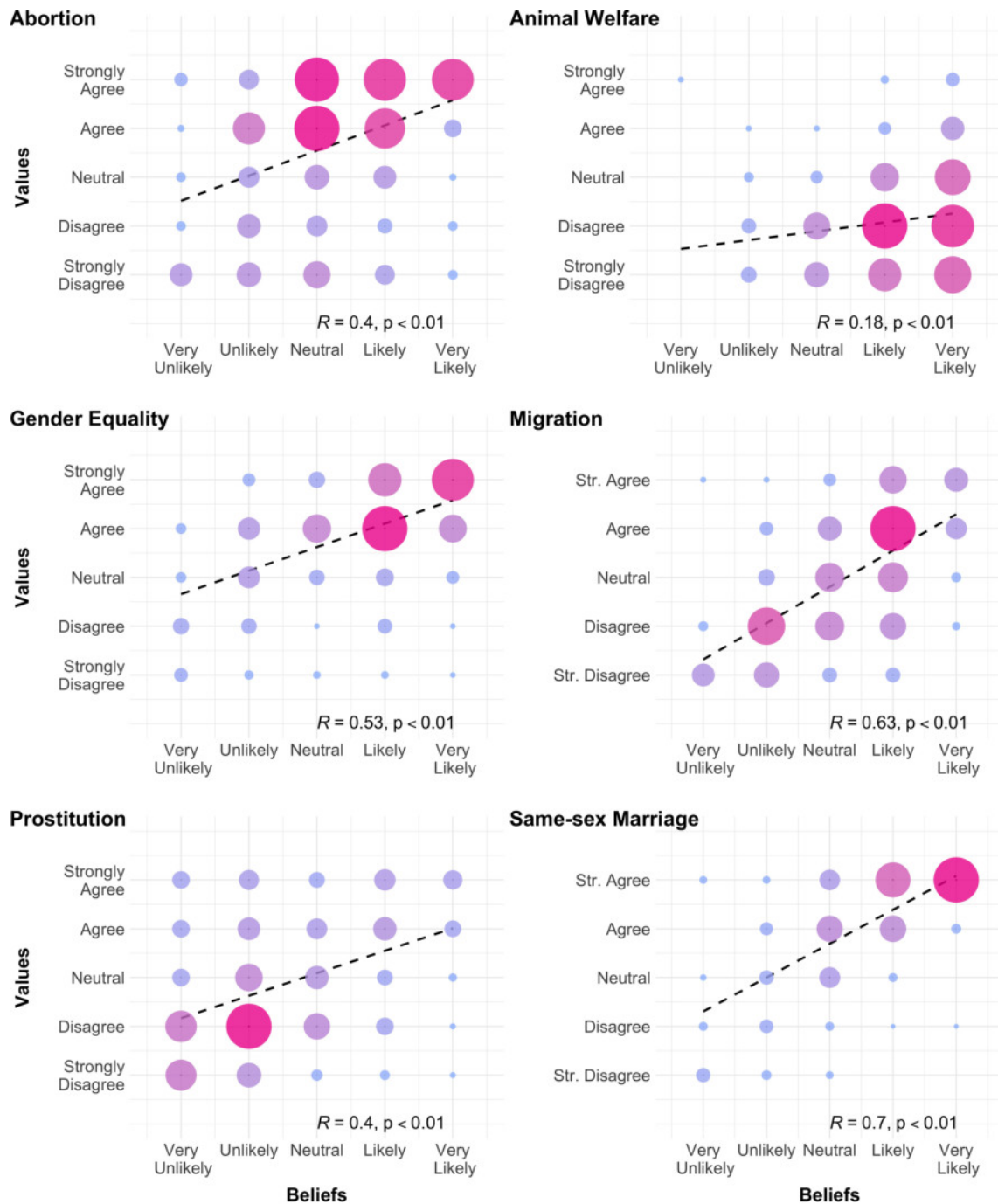
Table B1: Descriptive Statistics for Main Variables

	VALUENOT SALIENT	VALUE SALIENT	CONVINCE SELF	CONVINCE OTHER	BEING CONVINCED
Belief	3.46 (1.282)	3.42 (1.235)	3.48 (1.261)	3.53 (1.286)	3.44 (1.248)
Value		3.29 (1.398)	3.29 (1.425)	3.37 (1.441)	3.33 (1.436)
Donation	0.81 (1.019)	0.73 (0.998)	0.74 (1.007)		0.91 (1.073)
Pol. Attitude	4.76 (2.553)	4.71 (2.447)	4.52 (2.618)	4.42 (2.534)	4.58 (2.591)
<i>Political Affiliation (Share within treatment):</i>					
Democrats	0.45 (0.498)	0.47 (0.499)	0.43 (0.495)	0.45 (0.498)	0.49 (0.500)
Republicans	0.21 (0.404)	0.18 (0.382)	0.22 (0.413)	0.19 (0.393)	0.16 (0.368)
Independent	0.26 (0.436)	0.26 (0.440)	0.28 (0.449)	0.27 (0.446)	0.25 (0.434)
Other	0.03 (0.174)	0.02 (0.122)	0.02 (0.144)	0.03 (0.162)	0.04 (0.199)
None	0.06 (0.230)	0.08 (0.269)	0.05 (0.227)	0.05 (0.226)	0.05 (0.226)
Observations	2250	2310	2262	2178	2178

*Note:* The table shows the mean (share for political affiliations) and standard deviation in parentheses for each of the four treatments, VALUENOTSALIENT, VALUESALIENT, CONVINCESSELF, CONVINCEOther and the auxiliary treatment BEINGCONVINCED pooled over the different debates.

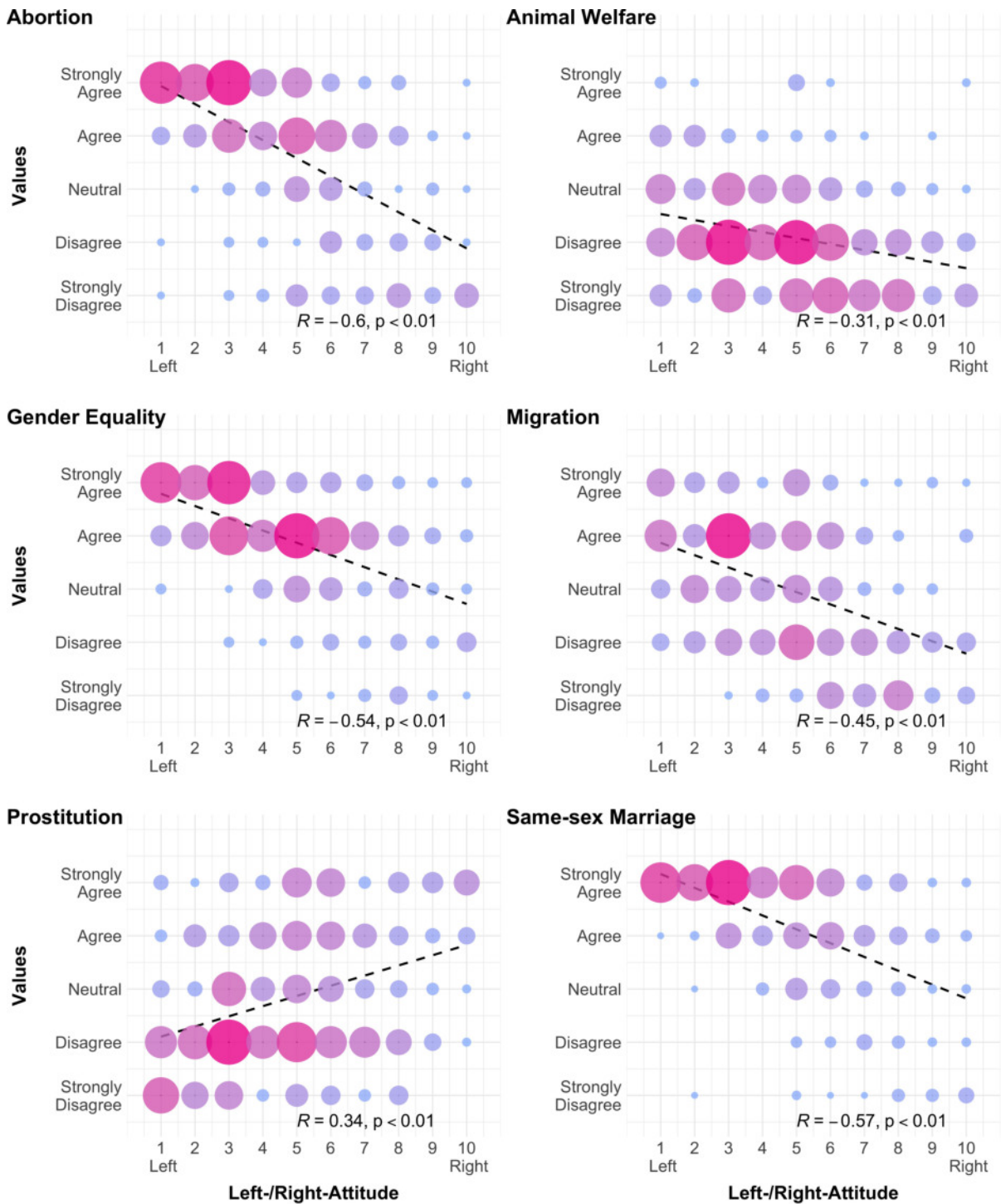
## B.2 Supplementary Results for Hypothesis 1

Figure B.4: Relationship between values and beliefs in VALUE SALIENT, by topic



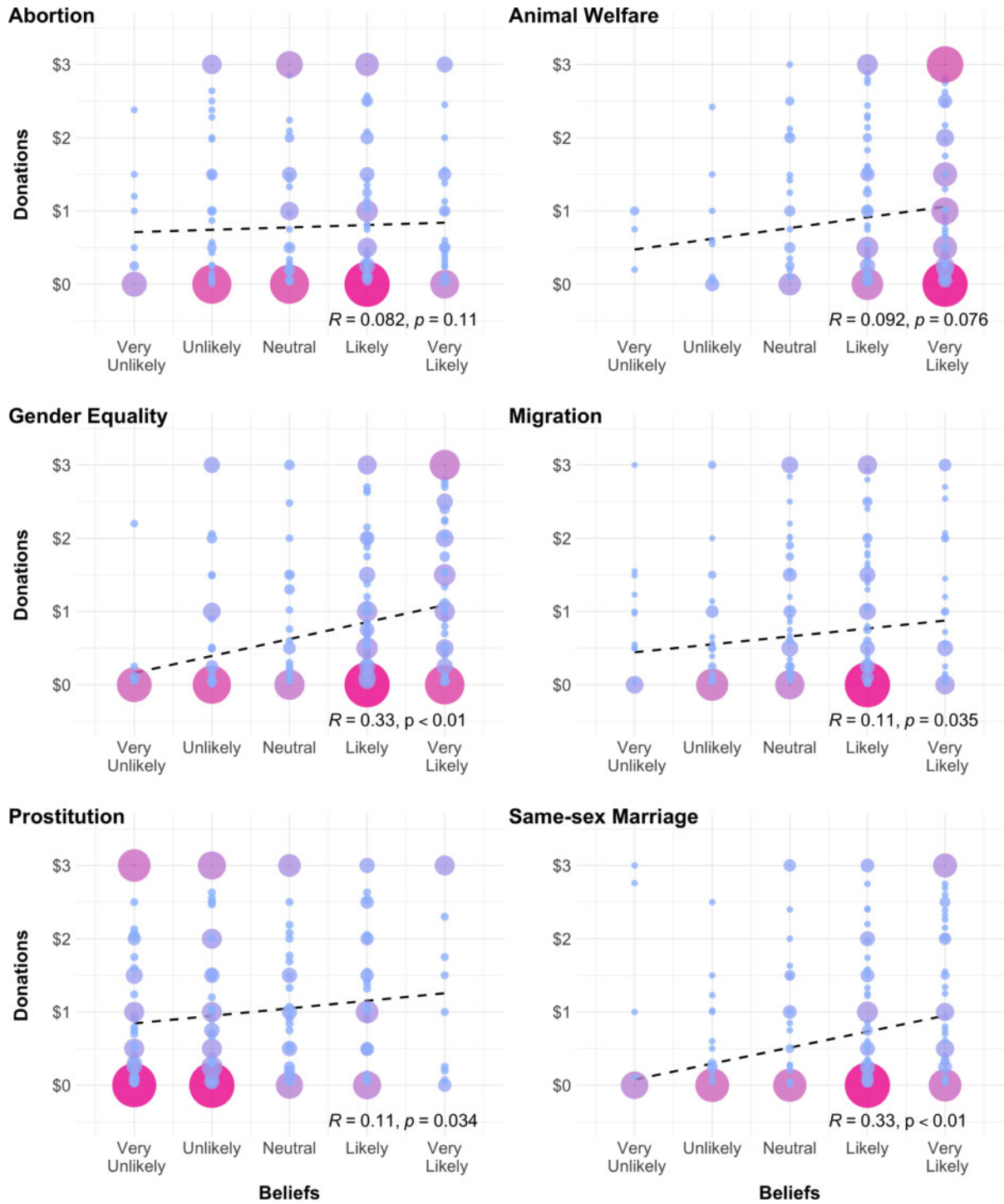
Note: Figure B.4 shows the correlation between values and beliefs in the VALUE SALIENT treatment, separately for each of the six policy debates. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of values on beliefs respectively political attitudes. The Spearman correlation coefficient, R, and its p-value are given at the bottom right of each graph.

Figure B.5: Relationship between political attitudes and beliefs in VALUE SALIENT, by topic



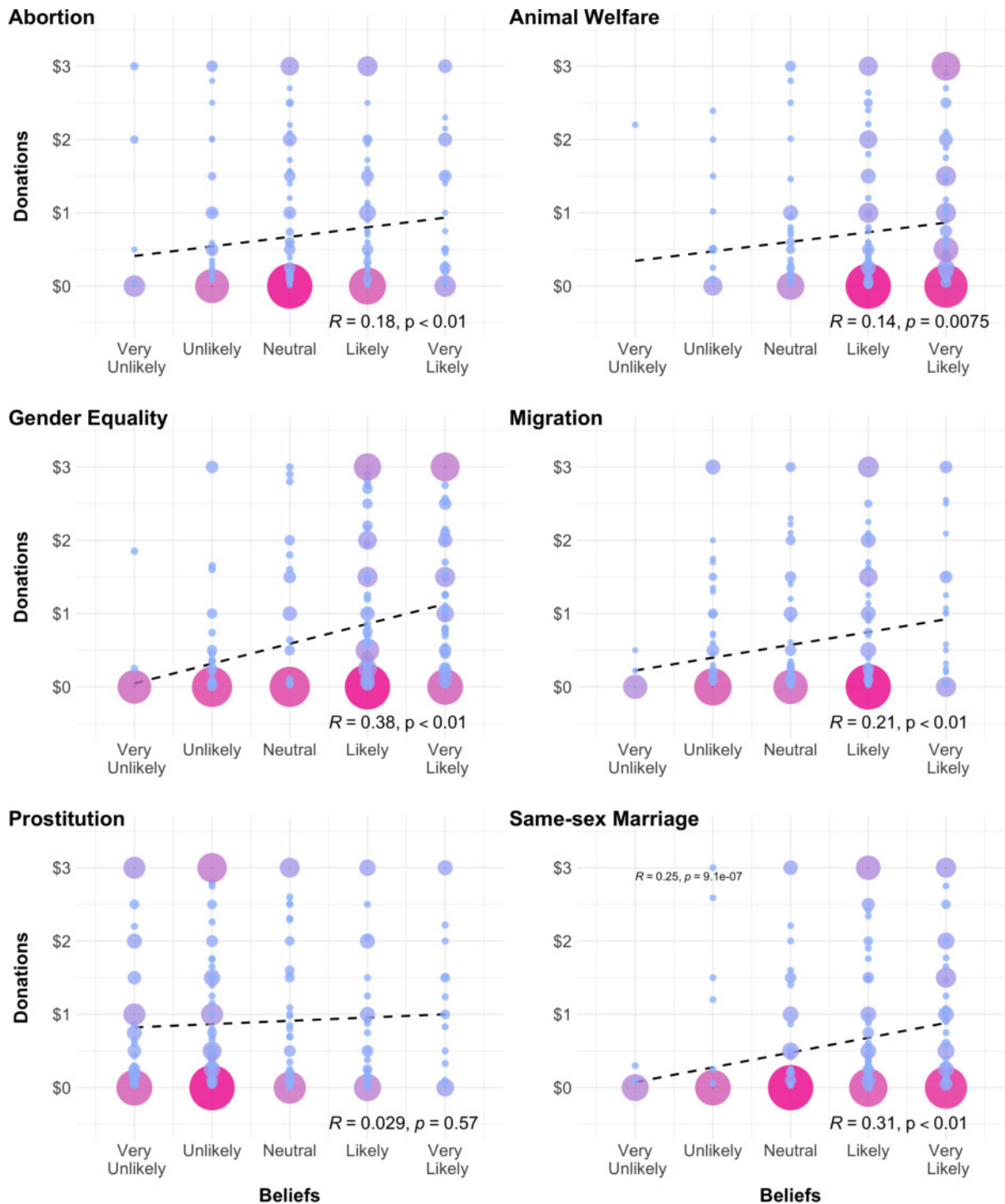
Note: Figure B.5 shows the correlation between moral values and political attitudes in the VALUE SALIENT treatment, separately for each of the six policy debates. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of values on beliefs respectively political attitudes. The Spearman correlation coefficient, R, and its p-value are given at the bottom right of each graph.

Figure B.6: Relationship between donations and beliefs in VALUENOTSALIENT, by topic



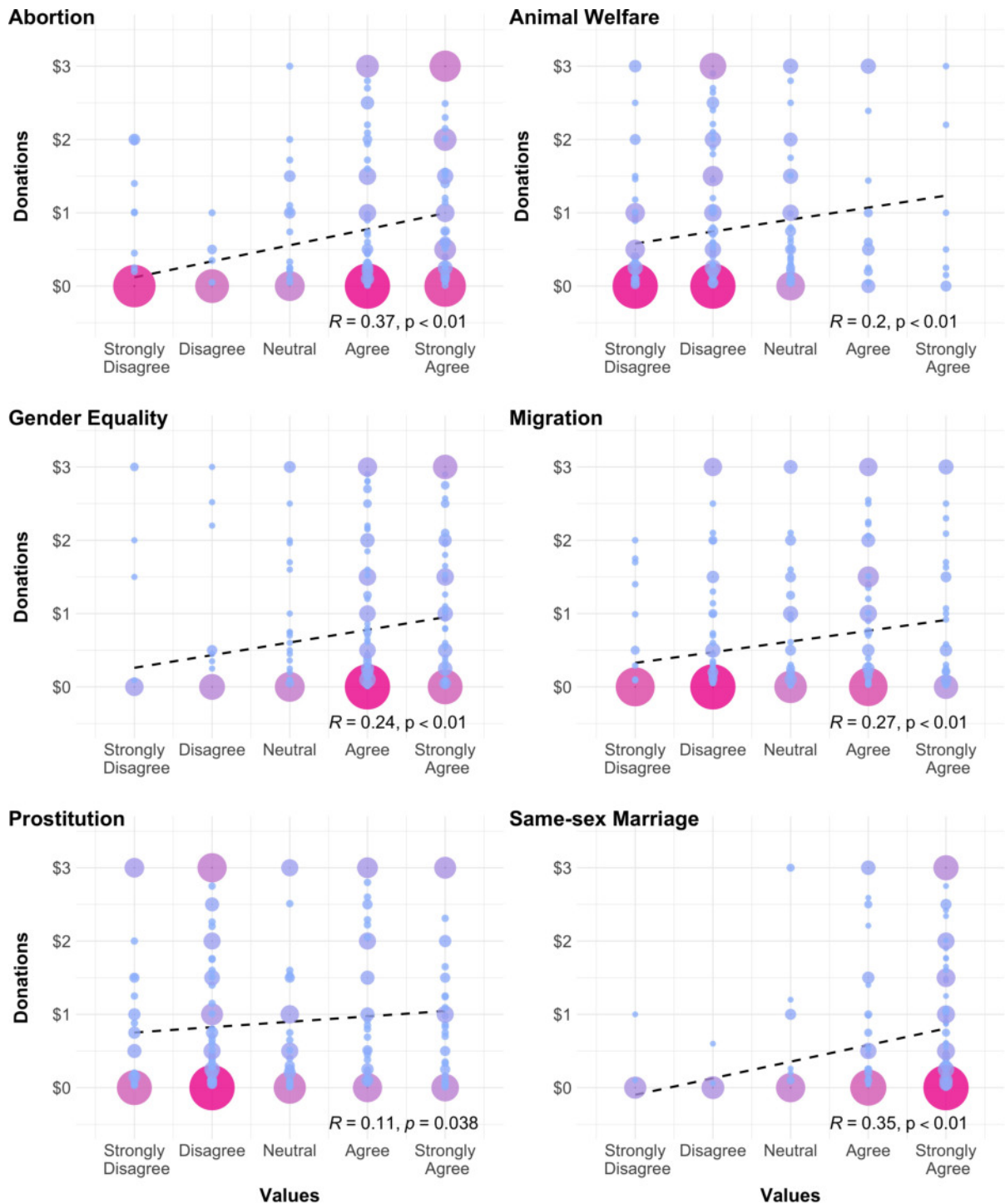
Note: The figure shows the correlation between beliefs and donations in treatment VALUENOTSALIENT separately for each of the six policy debates. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of donations on beliefs respectively values. The Spearman correlation coefficient, R, and its p-value are given at the bottom right of each graph.

Figure B.7: Relationship between donations and beliefs in VALUESALIENT, by topic



Note: The figure shows the correlation between beliefs and donations in treatment VALUESALIENT separately for each of the six policy debates. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of donations on beliefs respectively values. The Spearman correlation coefficient, R, and its p-value are given at the bottom right of each graph.

Figure B.8: Relationship between donations and values in VALUE SALIENT, by topic



Note: The figure shows the correlation between values and donations in treatment VALUE SALIENT separately for each of the six policy debates. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of donations on beliefs respectively values. The Spearman correlation coefficient, R, and its p-value are given at the bottom right of each graph.

### B.3 Supplementary Results for Hypothesis 2

Table B2: Robustness checks for Table 3

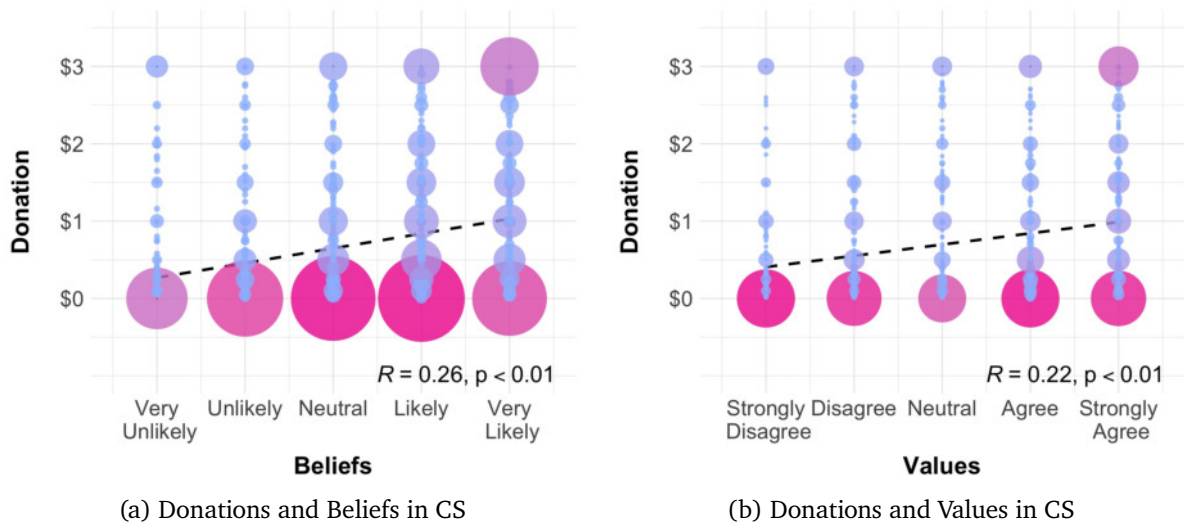
	(1)	(2)
VALUESALIENT	-0.104**	-0.162**
	[0.074]	[0.065]
Pol. Attitude ( $\tilde{p}$ )	0.301***	0.294***
	[0.062]	[0.054]
<b>VALUESALIENT</b>	<b>0.201***</b>	<b>0.258***</b>
<b>× Pol. Attitude</b>	<b>[0.092]</b>	<b>[0.083]</b>
Constant	3.285***	3.310***
	[0.050]	[0.0643]
Observations	2550	3006
Incl. Controls	No	No
Pol. Attitude ( $\tilde{p}$ )	Left-Right	Left-Right
	Scale	Scale
Variable	Left=1	Left=1

Standard errors clustered by individual in parentheses,  
 \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

*Note:* The table reports the results from running the regression specification from column (1) of Table 3, using the subsamples of column (2) and (3) from the same table. In column (1) we present the results using the subsample from column (2) of Table 3 and in column (2) we present the results using the subsample from column (1) of Table 3. Each regression uses the observations from treatments VALUENOTSALIENT (375 observations) and VALUESALIENT (385 observations), pooled over all six debates. VALUESALIENT is a dummy variable equal to one if the individual was assigned to treatment VALUESALIENT and hence equal to zero if the individual was assigned to treatment VALUENOTSALIENT.

## B.4 Supplementary Results on Beliefs, Values and Donations

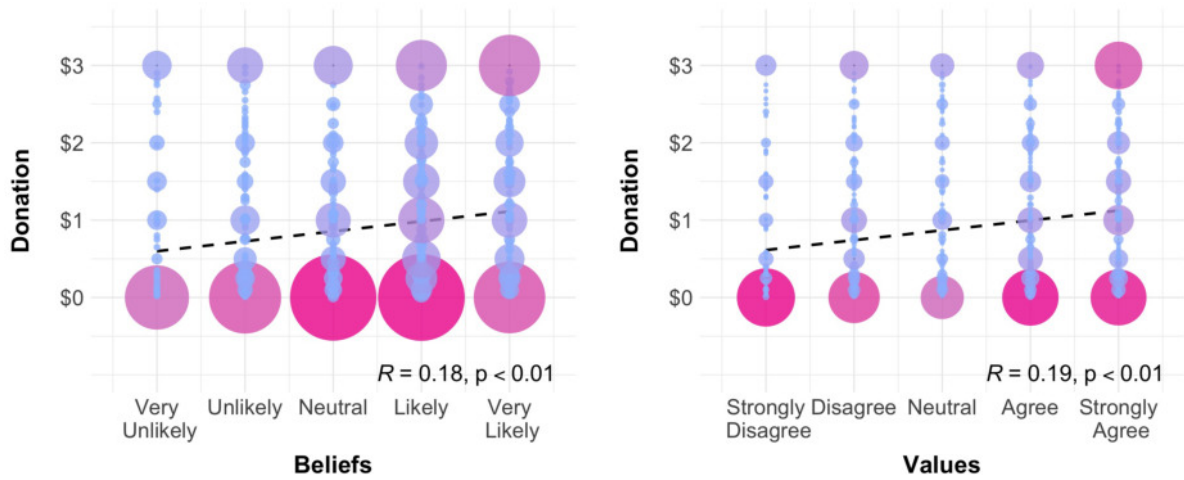
Figure B.9: Correlation between Donations and Values/Beliefs in Treatment CONVINCESSELF



Note: Figure B.9(a) shows the correlation between beliefs and donations in treatment CONVINCESSELF, Figure B.9(b) shows the correlation between values and donations in treatment CONVINCESSELF. The data points are weighted by the number of observations, which is reflected in both the color and the size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of donations on beliefs respectively values. The Spearman correlation coefficient,  $R$ , and its  $p$ -value are given at the bottom right of each graph.



Figure B.10: Correlation between Donations and Values/Beliefs in Treatment CONVINCETHER



(a) Donations and Beliefs in COR

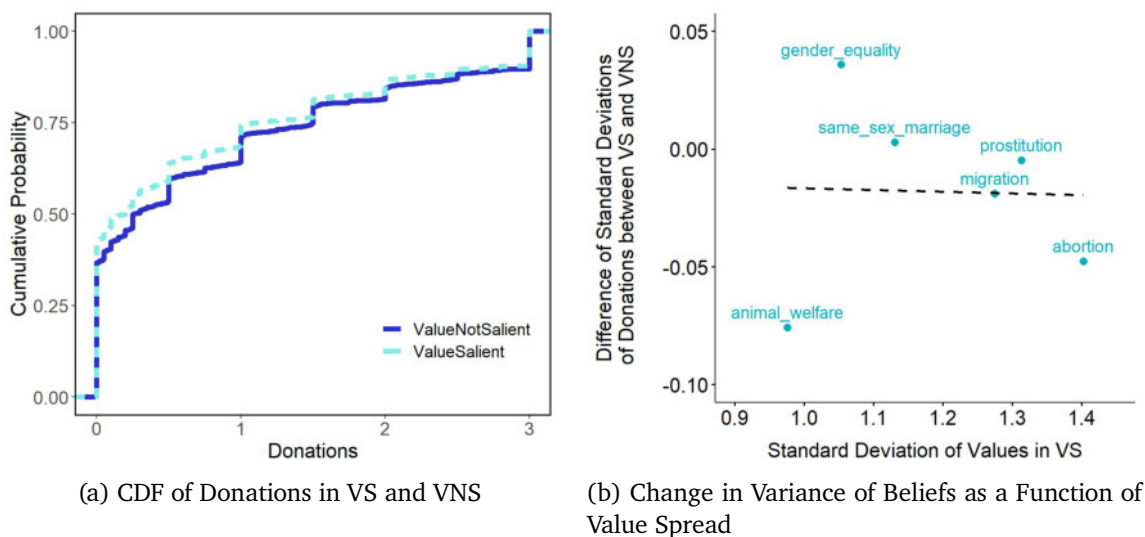
(b) Donations and Values in COR

Note: Figure B.10(a) shows the correlation between beliefs and donations in treatment CONVINCETHER, Figure B.10(b) shows the correlation between values and donations in treatment CONVINCETHER. The data points are weighted by the number of observations, which is reflected in both the color and size of the markers: the larger and redder the marker, the more observations; the smaller and bluer the marker, the fewer observations. The dotted line represents the result of a linear regression of donations on beliefs respectively values. The Spearman correlation coefficient,  $R$ , and its  $p$ -value are given at the bottom right of each graph.

### B.4.1 The Relationship between Hypothesis 2 and Donation Decisions

The results described in this section did not form part of our pre-registration. However, as an additional ex-post analysis, we provide documentation of the relationship between the donation decisions observed in the VALUE SALIENT and VALUE NOT SALIENT treatment conditions. The general conclusion of the results in this section is that donation decisions were not significantly impacted by the treatment variation. This indicates that although beliefs were shifted by the treatment, this shift did not translate into a change in donation behavior. This result, therefore, contributes to the growing literature that documents a complex relationship between measured beliefs and behavior. While some of the work in this literature documents evidence of beliefs causally affecting behavior in the manner predicted by standard economic models (see, e.g., Costa-Gomes, Huck, and Weizsäcker 2014; Haaland, Roth, and Wohlfart 2023), there is also body of work that show a divergence between predictions and behavior (see, e.g., Costa-Gomes and Weizsäcker 2008; Ivanov, Levin, and Niederle 2010; Haaland and Roth 2023).

Figure B.11: Results for Hypotheses 2a and 2b looking at donations

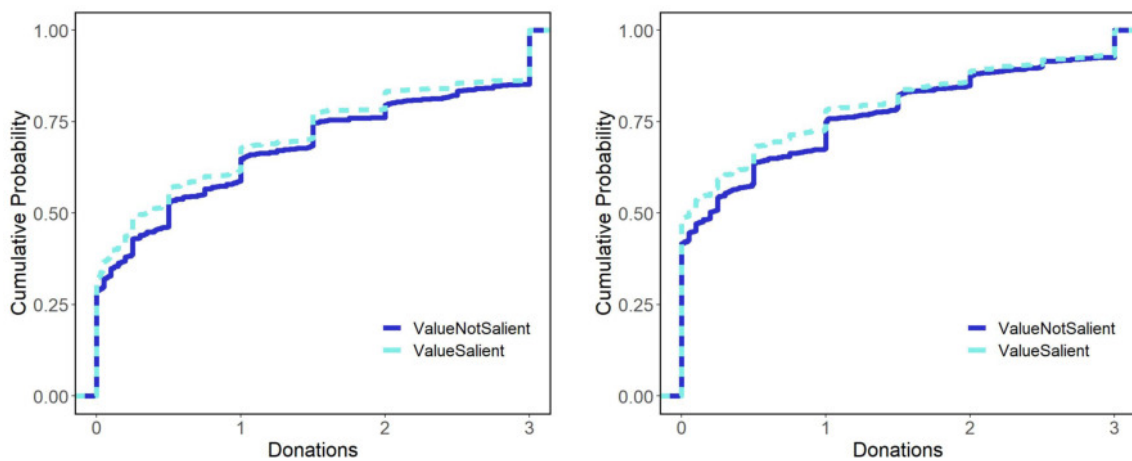


*Note:* Figure B.11 replicates the exercises in Figure 3, but replaces beliefs with donations. Figure B.11(a) shows the cumulative density function of donations in treatments VALUE SALIENT and VALUE NOT SALIENT. Figure B.11(b) shows the results from conducting the same exercise as we used to evaluate Hypothesis 2b, but examining donations instead of beliefs. The y-axis shows the difference of the standard deviations of donations between treatments VALUE SALIENT and VALUE NOT SALIENT and the x-axis shows the standard deviation of values in treatment VALUE SALIENT. The dotted line depicts the result from a linear regression of the difference of the standard deviations on the standard deviation of values.

In interpreting these results, it is important to keep in mind that we also observe a strong and robust correlation between donation decisions and both beliefs and values in each of our treatment conditions. There are several reasons why the shift in beliefs might not trans-

late directly into a shift in donation decisions, including the following. First, it may be the case that deep values are a more important driver of donation decisions than factual beliefs. This would explain the correlations between donations and beliefs and values (since values and beliefs are correlated), but would also be consistent with the fact that the shift in beliefs doesn't translate into a shift in donation decisions. Second, it is plausible that when individuals face their donation decision, the underlying contentious value debate is triggered and becomes salient at the point of making the donation decision. This would potentially negate the treatment differences generated by varying the salience of the value debates introduced by the VALUE-SALIENT and VALUE-NOT-SALIENT treatment conditions at the point of making the donation decision.

Figure B.12: Results for Hypotheses 2a and 2b looking at donations



(a) CDF of Donations in VS and VNS, subjects on the left of the political spectrum

(b) CDF of Donations in VS and VNS, subjects on the right of the political spectrum

*Note:* Figure B.12(a) shows the cumulative density function of donations in treatments VALUE-SALIENT and VALUE-NOT-SALIENT for individuals on the left of the political spectrum, i.e., below the mean of the political attitude variable. Figure B.12(b) shows the cumulative density function of donations in treatments VALUE-SALIENT and VALUE-NOT-SALIENT for individuals on the right of the political spectrum, i.e., above the mean of the political attitude variable.

Table B3: Influence of Increased Salience of Values on Donations

	(1)	(2)	(3)
VALUESALIENT	-0.071 [0.078]	-0.019 [0.110]	-0.151 [0.107]
Pol. Attitude ( $\tilde{p}$ )	0.284*** [0.090]	0.512*** [0.113]	0.306*** [0.107]
VALUESALIENT × Pol. Attitude	-0.013 [0.127]	-0.114 [0.156]	0.100 [0.150]
Constant	0.671*** [0.055]	0.453*** [0.080]	0.652*** [0.077]
Observations	4560	2550	3006
Pol. Attitude ( $\tilde{p}$ ) Variable	Left-Right Scale (Left = 1)	Party Affiliation (Democrat = 1)	Last Election (Clinton = 1)

Standard errors clustered by individual in parentheses, \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

*Note:* Each regression uses the observations from treatments VALUENOTSALIENT (375 observations) and VALUESALIENT (385 observations), pooled over all six debates. Smaller sample sizes in columns (2) and (3) result from missing information in the political affiliation variables and/or in the control variables. VALUESALIENT is a dummy variable equal to one if the individual was assigned to treatment VALUESALIENT and hence equal to zero if the individual was assigned to treatment VALUENOTSALIENT. We use three measures of the political attitudes variable. This is indicated in the last two rows of the table. In column (1), Political Attitude is a dummy equal to one if the individual is below the median on a 1 to 10 scale of political attitudes where 1 is the most left and 10 is the most right attitude. In column (2), Political Attitude is a dummy equal to one if the individual identifies as a Democrat rather than as a Republican and in column (3) Political Attitude equals one if the individual indicated that they voted for Clinton in the 2016 elections and zero if they voted for Trump.

# C Sample Balance

Table C1: Sample Balance

Variable	(1) VALUENOT- SALIENT Mean/(SE)	(2) VALUE- SALIENT Mean/(SE)	(3) CONVINCE- SELF Mean/(SE)	(4) CONVINCE- OTHER Mean/(SE)	(5) BEING- CONVINCED Mean/(SE)	T-test Diff. (1)-(2)	T-test Diff. (1)-(3)	T-test Diff. (1)-(4)	T-test Diff. (1)-(5)	T-test Diff. (2)-(3)	T-test Diff. (2)-(4)	T-test Diff. (2)-(5)	T-test Diff. (3)-(4)	T-test Diff. (3)-(5)	T-test Diff. (4)-(5)
Age	44.096 (0.821)	43.821 (0.798)	45.040 (0.809)	44.490 (0.828)	42.972 (0.830)	0.275	-0.944	-0.394	1.124	-1.219	-0.670	0.848	0.549	2.067*	1.518
Female	0.507 (0.026)	0.514 (0.026)	0.520 (0.026)	0.510 (0.026)	0.540 (0.026)	-0.008	-0.013	-0.003	-0.033	-0.006	0.005	-0.026	0.010	-0.020	-0.030
Ethnicity															
White	0.739 (0.023)	0.766 (0.022)	0.790 (0.021)	0.744 (0.023)	0.736 (0.023)	-0.028	-0.052*	-0.005	0.003	-0.024	0.022	0.031	0.047	0.055*	0.008
Asian	0.069 (0.013)	0.057 (0.012)	0.066 (0.013)	0.063 (0.013)	0.063 (0.013)	0.012	0.003	0.006	0.006	-0.009	-0.006	-0.006	0.003	0.003	0.000
Black	0.157 (0.019)	0.125 (0.017)	0.106 (0.016)	0.138 (0.018)	0.124 (0.017)	0.033	0.051**	0.020	0.033	0.019	-0.013	0.001	-0.032	-0.018	0.014
Mixed	0.019 (0.007)	0.034 (0.009)	0.021 (0.007)	0.025 (0.008)	0.052 (0.012)	-0.015	-0.003	-0.006	-0.034**	0.013	0.009	-0.019	-0.004	-0.031**	-0.028*
Other	0.016 (0.006)	0.018 (0.007)	0.016 (0.006)	0.030 (0.009)	0.025 (0.008)	-0.002	0.000	-0.014	-0.009	0.002	-0.012	-0.007	-0.014	-0.009	0.006
Employment Status															
Starting job next month	0.016 (0.006)	0.018 (0.007)	0.008 (0.005)	0.008 (0.005)	0.008 (0.005)	-0.002	0.008	0.008	0.008	0.010	0.010	0.010	-0.000	-0.000	0.000
Full-time	0.475 (0.026)	0.462 (0.025)	0.488 (0.026)	0.471 (0.026)	0.435 (0.026)	0.012	-0.013	0.004	0.039	-0.026	-0.009	0.027	0.017	0.053	0.036
Not in paid work	0.200 (0.021)	0.190 (0.020)	0.223 (0.021)	0.229 (0.022)	0.215 (0.022)	0.010	-0.023	-0.029	-0.015	-0.033	-0.039	-0.025	-0.006	0.008	0.014
Part-time	0.192 (0.020)	0.203 (0.021)	0.159 (0.019)	0.165 (0.020)	0.190 (0.021)	-0.011	0.033	0.027	0.002	0.043	0.037	0.013	-0.006	-0.031	-0.025
Unemployed	0.085 (0.014)	0.078 (0.014)	0.072 (0.013)	0.094 (0.015)	0.102 (0.016)	0.007	0.014	-0.008	-0.017	0.006	-0.016	-0.024	-0.022	-0.030	-0.008
Other	0.032	0.049	0.050	0.033	0.050	-0.017	-0.018	-0.001	-0.018	-0.001	0.016	-0.000	0.017	0.001	-0.017
Education															
Secondary education	0.005 (0.004)	0.013 (0.006)	0.024 (0.008)	0.028 (0.009)	0.017 (0.007)	-0.008	-0.019**	-0.022**	-0.011	-0.011	-0.015	-0.004	-0.004	0.007	0.011
High school diploma	0.232 (0.022)	0.182 (0.020)	0.218 (0.021)	0.209 (0.021)	0.256 (0.023)	0.050*	0.014	0.023	-0.024	-0.036	-0.028	-0.074**	0.008	-0.039	-0.047
Undergraduate degree	0.387 (0.025)	0.436 (0.025)	0.395 (0.025)	0.386 (0.026)	0.375 (0.025)	-0.050	-0.009	0.001	0.012	0.041	0.051	0.062*	0.010	0.021	0.011
Observations	375	385	377	363	363										

Continued on next page.

Table C1 – continued from previous page

Variable	(1) VALUENOT- SALIENT Mean/(SE)	(2) VALUE- SALIENT Mean/(SE)	(3) CONVINCE- SELF Mean/(SE)	(4) CONVINCE- OTHER Mean/(SE)	(5) BEING- CONVINCED Mean/(SE)	T-test Diff. (1)-(2)	T-test Diff. (1)-(3)	T-test Diff. (1)-(4)	T-test Diff. (1)-(5)	T-test Diff. (2)-(3)	T-test Diff. (2)-(4)	T-test Diff. (2)-(5)	T-test Diff. (3)-(4)	T-test Diff. (3)-(5)	T-test Diff. (4)-(5)
Technical/ comm. college	0.160 (0.019)	0.192 (0.020)	0.154 (0.019)	0.168 (0.020)	0.176 (0.020)	-0.032	0.006	-0.008	-0.016	0.038	0.024	0.016	-0.014	-0.022	-0.008
Graduate degree	0.181 (0.020)	0.143 (0.018)	0.172 (0.019)	0.165 (0.020)	0.132 (0.018)	0.038	0.009	0.016	0.049*	-0.030	-0.022	0.011	0.007	0.040	0.033
Doctorate degree	0.027 (0.008)	0.031 (0.009)	0.029 (0.009)	0.044 (0.011)	0.039 (0.010)	-0.005	-0.003	-0.017	-0.012	0.002	-0.013	-0.007	-0.015	-0.009	0.006
No formal qualification	0.005 (0.004)	0.003 (0.003)	0.003 (0.003)	0.000 (0.000)	0.003 (0.003)	0.003	0.003	0.005	0.003	-0.000	0.003	-0.000	0.003	-0.000	-0.003
Not applicable	0.003 (0.003)	0.000 (0.000)	0.005 (0.004)	0.000 (0.000)	0.003 (0.003)	0.003	-0.003	0.003	-0.000	-0.005	.n	-0.003	0.005	0.003	-0.003
Income															
Less than \$10000	0.067 (0.013)	0.052 (0.011)	0.066 (0.013)	0.066 (0.013)	0.061 (0.013)	0.015	0.000	0.001	0.006	-0.014	-0.014	-0.009	0.000	0.006	0.006
\$10000-\$15999	0.077 (0.014)	0.047 (0.011)	0.064 (0.013)	0.069 (0.013)	0.039 (0.010)	0.031*	0.014	0.008	0.039**	-0.017	-0.022	0.008	-0.005	0.025	0.030*
\$16000-\$19999	0.032 (0.009)	0.042 (0.010)	0.019 (0.007)	0.028 (0.009)	0.030 (0.009)	-0.010	0.013	0.004	0.002	0.023*	0.014	0.011	-0.009	-0.012	-0.003
\$20000-\$29999	0.117 (0.017)	0.094 (0.015)	0.111 (0.016)	0.099 (0.016)	0.129 (0.018)	0.024	0.006	0.018	-0.012	-0.018	-0.006	-0.036	0.012	-0.018	-0.030
\$30000-\$39999	0.107 (0.016)	0.119 (0.017)	0.106 (0.016)	0.085 (0.015)	0.096 (0.016)	-0.013	0.001	0.021	0.010	0.013	0.034	0.023	0.021	0.010	-0.011
\$40000-\$49999	0.107 (0.016)	0.078 (0.014)	0.093 (0.015)	0.124 (0.017)	0.110 (0.016)	0.029	0.014	-0.017	-0.004	-0.015	-0.046**	-0.032	-0.031	-0.017	0.014
\$50000-\$59999	0.091 (0.015)	0.101 (0.015)	0.095 (0.015)	0.102 (0.016)	0.102 (0.016)	-0.011	-0.005	-0.011	-0.011	0.006	-0.001	-0.001	-0.006	-0.006	0.000
\$60000-\$69999	0.056 (0.012)	0.078 (0.014)	0.088 (0.015)	0.066 (0.013)	0.058 (0.012)	-0.022	-0.032*	-0.010	-0.002	-0.010	0.012	0.020	0.021	0.030	0.008
\$70000-\$79999	0.085 (0.014)	0.086 (0.014)	0.077 (0.014)	0.088 (0.015)	0.069 (0.013)	-0.000	0.008	-0.003	0.016	0.009	-0.002	0.017	-0.011	0.008	0.019
\$80000-\$89999	0.045 (0.011)	0.052 (0.011)	0.042 (0.010)	0.052 (0.012)	0.050 (0.011)	-0.007	0.003	-0.007	-0.004	0.010	-0.000	0.002	-0.010	-0.007	0.003
\$90000-\$99999	0.040 (0.010)	0.047 (0.011)	0.040 (0.010)	0.052 (0.012)	0.047 (0.011)	-0.007	0.000	-0.012	-0.007	0.007	-0.006	-0.000	-0.013	-0.007	0.006
\$100000-\$149999	0.104 (0.016)	0.127 (0.017)	0.125 (0.017)	0.085 (0.015)	0.132 (0.018)	-0.023	-0.021	0.019	-0.028	0.003	0.042*	-0.005	0.039*	-0.008	-0.047**
More than \$150000	0.040 (0.010)	0.055 (0.012)	0.053 (0.012)	0.063 (0.013)	0.039 (0.010)	-0.015	-0.013	-0.023	0.001	0.001	-0.009	0.016	-0.010	0.014	0.025
Prefer not to say	0.032	0.023	0.021	0.019	0.039	0.009	0.011	0.013	-0.007	0.002	0.004	-0.015	0.002	-0.017	-0.019
Observations	375	385	377	363	363										

Continued on next page.

Table C1 – continued from previous page

Variable	(1) VALUENOT- SALIENT Mean/(SE)	(2) VALUE- SALIENT Mean/(SE)	(3) CONVINCE- SELF Mean/(SE)	(4) CONVINCE- OTHER Mean/(SE)	(5) BEING- CONVINCED Mean/(SE)	T-test Diff. (1)-(2)	T-test Diff. (1)-(3)	T-test Diff. (1)-(4)	T-test Diff. (1)-(5)	T-test Diff. (2)-(3)	T-test Diff. (2)-(4)	T-test Diff. (2)-(5)	T-test Diff. (3)-(4)	T-test Diff. (3)-(5)	T-test Diff. (4)-(5)
Prolific-Score	(0.009) 99.589 (0.064)	(0.008) 99.670 (0.059)	(0.007) 99.602 (0.066)	(0.007) 99.521 (0.088)	(0.010) 99.518 (0.075)	-0.081	-0.013	0.069	0.071	0.068	0.149	0.152	0.081	0.084	0.003
Left-right political attitude	4.763 (0.132)	4.706 (0.125)	4.523 (0.135)	4.419 (0.133)	4.579 (0.136)	0.056	0.240	0.344*	0.184	0.184	0.288	0.128	0.104	-0.056	-0.160
Observations	1863	375	385	377	363	363									

Note: The table shows the means and standard errors (in round parentheses) of demographic variables for the individuals in our sample for each of the 5 treatments as well as pairwise comparisons of the means using a t-test.

## D Preregistration Document

### Morals, Beliefs, and Actions

Kai Barron (WZB), Anna Becker (UCL), Steffen Huck (UCL and WZB)

*This version: 14/01/2020*

#### PART I: EXPERIMENTAL DESIGN

##### Setup of Experiment

This experiment will be run as an online survey. The sample for the four main treatments will consist of 1,500 individuals that are representative of the US population in terms of age, sex, and ethnicity. An additional 375 subjects will be recruited later for an auxiliary treatment, Treatment 4b, which builds on subjects' choices in Treatment 4a. All participants will be recruited via the online platform Prolific. Subjects will be paid £3 for participation and have the option to earn a bonus in Treatments 1, 2, 3 and 4b.<sup>28</sup> A strict no-deception policy will be followed. The experiment is programmed using the experimental software o-Tree (Chen, Schonger, and Wickens (2016)).

##### Experimental Design

At the beginning of the experiment, subjects are randomized into one out of four treatment groups such that each group consists of 375 subjects. The four treatments are described in the following text.

###### *Treatment 1 – “Control”*

###### *Part 1*

Subjects are asked to state how likely they think a statement that they are presented with is true. The statements have been chosen such that they can be associated with a policy domain and typically refer to facts on which scientific consensus has not been reached yet. Subjects use a five-point Likert scale to indicate their beliefs. They can choose between the following options: “Very Unlikely”, “Very Likely”, “Neutral”, “Likely” and “Very Likely”. After a waiting time of 15 seconds subjects can proceed to the next page.

This is repeated six times for six different policy domains. These are migration, animal welfare, gender equality, abortion, prostitution and gay rights. Table C1 in the Appendix provides an overview over all domains and the statements presented to subjects. The order with which the statements on the different domains are shown to participants is randomised at the individual level.

---

28. The show-up fee is converted into US dollars, since the subjects are recruited from the USA.



## *Part 2*

After subjects have submitted their beliefs on the six different issues, they are informed that they have the option to make six donations to charities. They are also informed that one out of the six decisions they are about to make will be chosen at random to be implemented.

Subjects then see the statement they were earlier confronted with, the belief that they stated and the option to donate to a charity that is active in the respective policy domain. As in Part 1, this process is repeated six times following the same order of domains as in Part 1. Subjects are informed briefly about the aims of the charities and can indicate with a slider how much they would like to donate. They are provided with \$3 and can choose to donate any amount between \$0 and \$3. Subjects will be paid the remaining amount (i.e. what they decided not to donate) at the end of the experiment. They can proceed to the next page at any time after a waiting time of 15 seconds.

## *Treatment 2 – “Moral Values”*

### *Part 1*

Treatment 2 is similar to Treatment 1, with the following exceptions.

Different to Treatment 1, in Treatment 2 subjects will also be asked to state some moral values that are related to the same six policy domains. The question of how much they agree or disagree with a moral statement they are presented with appears above the question regarding the factual statement. Subjects use a five-point Likert scale to indicate their agreement. They can choose between the following options: “Strongly Disagree”, “Disagree”, “Neutral”, “Agree” and “Strongly Agree”. The moral statements can be found in Table C1 in the Appendix.

### *Part 2*

Different to Treatment 1, in Treatment 2, subjects are also reminded of the moral values they stated in Part 1. The screen will therefore show the moral statement and the subject’s choice above the factual statement and the subject’s decision and then offers the subjects to donate to a proposed charity.

## *Treatment 3 – “Convincing Yourself”*

As in Treatment 2, subjects are asked to state moral values, factual beliefs, and then make a charitable donation. The key difference between Treatment 3 and Treatment 2 is that in Treatment 3 the moral statement and the factual statement are presented to subjects on the same screen as the option to donate to charity, which is presented at the bottom of the page. They receive the same information as subjects in Part 2 of Treatments 1 and 2 and

face the same charitable giving decision, but in Treatment 3 all three decisions are made on the same page (i.e. the beliefs, moral value judgments, and charitable donations).

As in the other treatments, this will be repeated six times for the six different policy domains where the order is randomized on the individual level. Subjects use a Likert scale to indicate their values and beliefs and a slider to indicate how much they would like to donate. They are provided with \$3 per decision. One of those decisions will be chosen at random to be implemented of which subjects are informed about in advance. They will be paid what they decide not to donate after the experiment.

#### *Treatment 4a – “Convincing Others”*

As in Treatment 2, subjects are asked to state their moral values and their factual beliefs. The moral statement and the factual statement are presented to subjects on the same screen underneath each other as in Part 1 of Treatment 2.

Before stating values and beliefs, the subject is informed that another participant will have the option to donate to a related charity. Importantly, the other participant will make their charitable donation decision after being informed about the moral values and factual beliefs that the subject reports (i.e. the moral values and factual beliefs reported by subjects in Treatment 4a will be sent to subjects in Treatment 4b before participants in 4b make their charitable donation decisions).

Subjects in Treatment 4a are presented with the information the other participant in Treatment 4b will be shown about the charity. The donation decision that will be completed by Treatment 4b subjects is the same as in the previous Treatments, i.e. the participant has \$3 available of which they keep what they decide not to donate. Both the subjects in Treatment 4a and the other participant in Treatment 4b will be informed in advance that one of the six decisions will be chosen at random to be implemented.

*After the main Treatments 1, 2, 3 and 4a have been run, another 375 subjects will be recruited for Treatment 4b:*

#### *Treatment 4b – “Being Convinced”*

##### *Part 1*

This will be identical to Part 1 of Treatment 2.

##### *Part 2*

In Treatment 4b, Part 2 will be similar to Treatment 2, with the exception that instead of subjects being reminded of their own decisions regarding the moral and the factual statements,

subjects are now informed about the decisions of a participant from Treatment 4a when they make their charitable donation decision. As before, subjects will be provided with \$3 for each decision of which they will be paid what they decide not to donate. They are also informed that one out of their six decisions will be chosen at random to be implemented.

### **Post-experimental Survey**

After the experiment, all subjects will be asked to fill out a survey. The survey covers the following topics:

1. Personal Details
2. Political Attitude
3. Religious Attitude
4. Moral Foundations Questionnaire<sup>29</sup>
5. Questions on Moral Behavior
6. Cognitive Reflection Test

## **PART II: ANALYSIS PLAN**

### **II.1) Introduction**

Standard theories on belief formation typically disregard the desire of individuals to gather, avoid or interpret information in a way that serves non-instrumental purposes. A large recent literature has, however, shown that, for example, self-serving biases, wishful thinking and motivated reasoning are important determinants of belief formation. This project studies individuals' moral values as a potential source for motivated cognition and links it to partisan disagreement about factual statements, i.e. polarisation.<sup>30</sup>

We proceed in three steps. In each step, we test a small set of hypotheses that are inter-linked by a common underlying idea. In the first step, we seek to test whether individuals report moral values and factual beliefs that are aligned in the domain of political issues. Potential reasons why individuals might do this include: i) avoiding emotional discomfort, or cognitive dissonance emanating from holding or stating incoherent values and beliefs, and ii) using value and belief statements to justify self-interested actions (e.g. actions that increase the individual's material wealth). From the analysts' perspective, the presence of such belief-value constellations would provide a basis for taking an individual's moral values into consideration in trying to understand belief formation regarding factual statements. A

---

29. <https://moralfoundations.org/questionnaires/>, accessed 31/12/2019.

30. See for example Rabin (1995) for a theory on self-serving biases in moral reasoning.

potential motive for this could be a desire to establish something akin to a “moral identity” (Bénabou and Tirole (2011)).

In the second step, we then test whether there is a systematic pattern in the way factual beliefs are formed that may be partially responsible for the creation of these belief-value constellations (i.e. we ask whether factual beliefs are constructed in a way that forms these belief-value constellations). To do this, we compare the distribution of beliefs of individuals that were previously asked about, and hence reminded of, their related values (Treatment 2) with the distribution of beliefs of individuals in the control Treatment 1 (where no moral value statements are reported prior to stating factual beliefs). This allows us to assess the influence of being primed to think about the particular factual belief in question through the lens of the related value debate. Following on from this, we ask whether this mechanism can explain the recent trend of a polarization of beliefs in society that has been demonstrated to run along ideological lines (see, e.g., Gentzkow (2016)). In particular, there appears to be an increased disagreement about objective facts among members of society that is associated with political attitudes.<sup>31</sup> We hypothesize that the heterogeneity in moral values between different political groups may be leading to the formation of these polarized factual beliefs. Hence, we test whether factual beliefs become polarized as a function of political attitudes (e.g. as a result of individuals’ desire to adjust their beliefs to their moral values.)<sup>32</sup>

Lastly, we study two potential forces that might increase or decrease the degree of polarization of (stated) beliefs. First, we look at the impact of financial incentives (through motivated reasoning or self-persuasion), and second we will study the role of persuading others, which is particularly relevant in political contexts.

The first channel which considers the role of financial incentives on belief formation is important because there are many reasons why individuals might face costs to hold certain beliefs or values. For example, it may be costly to hold different beliefs and values to those held by individuals in one’s peer network.<sup>33</sup> Alternatively, holding particular beliefs and values may be costly when they induce the individual to take a particular costly action. For example, an individual who advocates the merits reducing inequality in society may feel compelled to take actions that reduce their own wealth in order to increase the wealth of a poorer individual. Rather than studying abstract costs (e.g. incoherence with one’s peers’

---

31. The most prominent current example is probably that of climate change where there is a widening gap in the views on the scientific evidence between Republicans and Democrats in the US (see e.g. McCright and Dunlap (2011) ).

32. In his theoretical work, Le Yaouanq (2023) links heterogeneity in political attitudes to partisan disagreement about objective facts through people’s idiosyncratic preferences regarding the policy implications of scientific findings. Our work seeks to understand the underlying psychological mechanisms in more detail.

33. On the role of group identity in belief polarization see e.g. Gennaioli and Tabellini (2018).

beliefs), we focus on the latter type of costs that accrue due to taking actions that reduce one's personal payment. In particular, we consider costly donation decisions and hypothesize that this cost leads individuals to bias their (stated) beliefs and values when they expect a related donation decision, with the bias operating in the opposite direction to the beliefs and values consistent with a higher donation.<sup>34</sup>

Second, we look at the potential role of individuals' desire to convince others to take actions that are in line with their own values. This motive of convincing others might lead people to further align their (stated) beliefs with their political agenda or goals, and perhaps to exaggerate these stated beliefs. We study whether subjects adjust their beliefs to be more extreme in order to convince another participant to give more or less to a suggested charity. In this case, we test whether individuals overstate the strength of their beliefs when trying to persuade another person to act in a certain way.

Section II.2) introduces the necessary notation, before we formalize our hypotheses in Section II.3).

## D.1 Notation

Let  $b_t$  denote the factual beliefs stated by individuals in Treatment  $t$ , where  $t \in \{1, 2, 3, 4a, 4b\}$ ,  $v_t$  are the moral values stated by individuals in treatment  $t$  where  $t \in \{2, 3, 4a, 4b\}$  and  $d_t$  are the donation decisions of individuals in treatment  $t$  where  $t \in \{1, 2, 3, 4b\}$ . Let  $F_{b_t}$  denote the cumulative distribution function (cdf) of factual beliefs in Treatment  $t$  where  $t \in \{1, 2, 3, 4a, 4b\}$ ,  $F_{v_t}$  the cdf of moral values in Treatment  $t$  where  $t \in \{2, 3, 4a, 4b\}$ , and  $F_{d_t}$  the cdf of donations in Treatment  $t$  where  $t \in \{1, 2, 3, 4b\}$ .

Let  $p_t$  denote the left-right political stance of individuals in Treatment  $t$  which will be elicited for all participants in the post-experimental survey using a Likert scale ranging from 1 to 10 where 1 is left and 10 is right, i.e.  $p_t$  is increasing in the degree to which an individual positions herself on the right of the political spectrum.  $F_{p_t}$  denotes the respective cdf. Belief and value statements were chosen such that the factual statement being true would provide support for agreement to the value statement. All statements are coded this way for the analysis but not necessarily presented to participants like this (Table C1 in the Appendix shows how statements are presented to subjects during the experiment).

At the same time, we recode all the moral value variables such that they are likely to be increasingly appealing as one moves from the right to the left of the political spectrum.

---

34. It is also possible that there are individuals that now bias their values and beliefs in the same direction as beliefs and values consistent with a higher donation as they consider them to be more important when relevant to justify a charitable donation. Our prior, however, is that the effect described above dominates.

Similarly, we code the factual belief variables such that if they are true, they support moral value positions typically held by individuals on the political left.<sup>35</sup> Charities are chosen such that if the moral statement is supported it would justify the charities' objectives.<sup>36</sup>

## D.2 Hypotheses

### D.2.1 Belief-Value Constellations

As stated above, we begin by testing whether individuals report moral values and factual beliefs on political issues that are aligned, whether their moral values are aligned with their political attitudes and whether factual beliefs and moral values are related to decision making in the form of costly charitable donation choices.

#### Hypothesis 1

- a) Moral values are positively correlated with beliefs:

$$\text{Corr}(v_2, b_2) \geq 0.$$

- b) Moral values are negatively correlated with political attitudes:

$$\text{Corr}(v_2, p_2) \leq 0.$$

- c) Donations are positively correlated with beliefs and values:

$$\text{Corr}(d_1, b_1) \geq 0, \text{Corr}(d_2, b_2) \geq 0, \text{Corr}(d_2, v_2) \geq 0.$$

Part a) of Hypothesis 1 tests whether moral values and factual beliefs reported by subjects in Treatment 2 are aligned. Recall that in Treatment 2 subjects are not aware of the opportunity to donate by the time they state their values and beliefs.

Part b) tests for a negative correlation between moral values and political attitudes. For example, Enke, Rodriguez-Padilla, and Zimmermann (2016) show that an individuals' moral type is strongly correlated with their political affiliation. Rather than looking at predefined moral types we look at concrete moral convictions with regard to certain policy domains.

---

35. The policy domain "Prostitution" is an ambiguous case. Individuals from the political left could both support and reject more liberal prostitution rights.

36. Therefore, it is worth pointing out that variables are coded such that a higher value of  $b_t$ ,  $v_t$ , and  $d_t$  should be consistent with a lower value of  $p_t$  according to the researcher team's priors.

We expect that the further an individual is on the left of the political spectrum (i.e. the lower is  $p_2$ ) the more likely they are to agree with the moral value statement.

In Part c) of Hypothesis 1, we hypothesize that individuals donate more when their moral values and their beliefs are such that the cause of the charity is justified by them (i.e. that donations are positively correlated with beliefs and moral values consistent with the charity's mandate).

### D.2.2 Construction of Beliefs

The following hypothesis tests: (i) whether the formation of factual beliefs is influenced by the values individuals hold, thereby explaining the formation of belief-value constellations, and (ii) whether this influence of values on belief formation may lead to a polarization of factual beliefs across the political spectrum.

#### Hypothesis 2

- a) The distribution of factual beliefs in Treatment 1 is different from the distribution of factual beliefs in Treatment 2:

$$F_{b_1} \neq F_{b_2}.$$

- b) Comparing across the six domains indexed by  $m$ , the difference between the variance in beliefs in Treatment 2 and the variance in beliefs in Treatment 1 is increasing in the variance in values in Treatment 2:

$$\frac{d[\text{Var}(b_2^m) - \text{Var}(b_1^m)]}{d[\text{Var}(v_2^m)]} \geq 0.$$

- c) Beliefs in Treatment 2 are more polarized than beliefs in Treatment 1:

$$E(b_2|p_2 < E(p_2)) - E(b_1|p_1 < E(p_1)) \geq E(b_2|p_2 > E(p_2)) - E(b_1|p_1 > E(p_1)).$$

Part a) of Hypothesis 2 exploits the fact that unlike subjects in Treatment 2, subjects in Treatment 1 are not asked to state their moral values. We test whether reminding subjects of their moral values has an impact on the distribution of factual beliefs which would indicate that moral values are relevant for the construction of beliefs. In Part b) of Hypothesis 2, we go further and test whether values exert a systematic influence on belief formation. In particular, Part b) of Hypothesis 2 posits that when there is more dispersion in the values that subjects hold with regard to a certain policy domain, we expect that there will be a shift towards more extreme beliefs as subjects are drawn towards more coherent belief-value

constructions.

The last part of Hypothesis 2 tests whether a polarization in beliefs can be explained as a result of the posited impact of values on belief formation. The inequality in Part c) states that the difference in the means of beliefs between Treatment 2 and Treatment 1 is greater for subjects below the mean of political attitudes, i.e. relatively on the left of the political spectrum, than for subjects above the mean of political attitudes, i.e. relatively on the right of the political spectrum. To put this another way, the hypothesis states that individuals on the left will increase their beliefs between Treatment 1 and Treatment 2 more than individuals on the right of the political spectrum, on average. Note that this also includes the case where subjects adjust their beliefs downwards (i.e. it is completely consistent with individuals on the right shifting their beliefs downward between Treatment 1 and 2, which would make the right-hand side negative).

If political attitudes are sufficiently widely dispersed, we would expect a positive left-hand side and a negative right-hand side, i.e. what we traditionally refer to as polarization. Otherwise, we expect to see a mild form of polarization where beliefs are adjusted to different extents by those on the left and the right of the political spectrum within our sample.

### D.2.3 Convincing Yourself and Convincing Others

The last hypothesis is split in two parts. In Part A, we look the role of self-serving biases that allow subjects to justify selfish behaviour and are expected to lead to a downward<sup>37</sup> bias in beliefs, values and charitable donations. Part B, on the other hand, studies whether introducing the opportunity to convince another participant to take an action that is in line with one's moral values can lead to a greater polarization of beliefs.

## Hypothesis 3

### A. Convincing Yourself

a) When there is an increase in the cost of holding certain beliefs and values individuals may shift their stated beliefs and values in the opposite direction:

i)  $b_2$  first-order stochastically dominates  $b_3$ , i.e.  $F_{b_2} \leq F_{b_3}$ .

ii)  $v_2$  first-order stochastically dominates  $v_3$ , i.e.  $F_{v_2} \leq F_{v_3}$ .

---

37. Here, "downward" refers to a bias in the direction that is consistent with being self-serving. In terms of the way we have defined our variables, it will also refer to lower values of  $b_t$ ,  $v_t$ , and  $d_t$ .



b) Donations in Treatment 3 are lower than in Treatment 2:

$$E(d_2) \geq E(d_3).$$

## B. Convincing Others

c) Beliefs in Treatment 4a are polarized in comparison to beliefs in Treatment 2:

$$E(b_{4a}|p_{4a} < E(p_{4a})) - E(b_2|p_2 < E(p_2)) \geq E(b_{4a}|p_{4a} > E(p_{4a})) - E(b_2|p_2 > E(p_2)).$$

In Treatments 1, 2 and 3 subjects are given the opportunity to donate to charity. Donating comes at the cost of a foregone bonus payment. Previous work has shown that people develop self-serving biases in order to excuse their selfishness in charitable giving (see e.g. Exley (2015) on the role of risk or Exley (2020) on using charity performance metrics as an excuse). In Part a) of Hypothesis 3 we test whether subjects shift their values and beliefs downwards to justify smaller donations in order to receive higher bonuses.

In Part b), we hypothesize that individuals donate less on average in Treatment 3 than in Treatment 2. In both treatments, a higher donation reduces the payment that the subject receives herself (i.e. the material incentives are identical across the two treatments). However, in Treatment 2, individuals only learn about the possibility to donate after they have already stated their values and beliefs which rules out the opportunity to make any adjustments to one's values in order to justify self-interested charitable donation decisions.<sup>38</sup>

The last part of Hypothesis 3 (i.e. Part c) tests whether individuals report more polarized beliefs when they have the opportunity to convince someone to make a donation which is the case in Treatment 4a. As in Hypothesis 2 c), the inequality in Hypothesis 3c) states that the difference in the means of beliefs between Treatment 4a and Treatment 2 is greater for subjects below the mean of political attitudes, i.e. relatively on the left of the political spectrum, than for subjects above the mean of the political spectrum, i.e. relatively on the right of the political spectrum.

As a spillover from Hypothesis 3 and Treatment 4a we can also study the effect of persuasion

---

38. There is the possibility that there exists a subset of individuals in Treatment 3, who instead of shifting down their reported beliefs and values to justify a lower charitable donation decision, instead shift up their values and beliefs in order to then enhance the signaling value of their donations. These individuals might then also donate more when facing Treatment 3 in comparison to Treatment 2. This effect would operate in the opposite direction to the main effect hypothesized in Hypothesis 3.b). For simplicity, we have stated Hypotheses 3.b) in terms of the average effect for the entire sample, working under the assumption that the main hypothesized effect of adjusting one's beliefs and values downwards in a self-serving fashion will dominate this potential countervailing secondary effect.

on recipients of the persuasion messages by looking at Treatment 4b. More specifically, we can study whether donations to charity increase when an individual sees their own values and beliefs confirmed by another person. We would expect that there will be a polarization of donation decisions when the political attitudes of the sender and the receiver are aligned.<sup>39</sup> In cases where the sender's and the receiver's political attitude are not aligned, we can think of (at least) two opposing possible effects. Individuals might either doubt their own convictions and reassess them or they might want to prove the sender wrong by exaggerating or lowering the donated amount. Ex ante, it is unclear which effect is expected to dominate. This in turn, might also depend on political affiliation. We do not propose any hypothesis here, but will document the data in an exploratory analysis.

---

39. This could be tested similarly to Part c) of Hypotheses 2 and 3.

## References

- Bénabou, Roland, and Jean Tirole. 2011. "Identity, Morals, and Taboos: Beliefs as Assets." *The Quarterly Journal of Economics* 126 (2): 805–855.
- Chen, Daniel L., Martin Schonger, and Chris Wickens. 2016. "oTree – An Open-source Platform for Laboratory, Online, and Field Experiments." *Journal of Behavioral and Experimental Finance* 9:88–97.
- Enke, Benjamin, Ricardo Rodriguez-Padilla, and Florian Zimmermann. 2016. "Moral Universalism and the Structure of Ideology." *Working Paper*.
- Exley, Christine L. 2015. "Excusing Selfishness in Charitable Giving: The Role of Risk." *Review of Economic Studies* 83, no. 2 (October): 587–628.
- . 2020. "Using Charity Performance Metrics as an Excuse Not to Give." *Management Science* 66 (2): 553–563.
- Gennaioli, Nicola, and Guido Tabellini. 2018. *Identity, Beliefs, and Political Conflict*. Technical report. Working Paper.
- Gentzkow, Matthew. 2016. *Polarization in 2016*. Technical report. Toulouse Network for Information Technology Working Paper.
- Le Yaouanq, Yves. 2023. "A model of voting with motivated beliefs." *Journal of Economic Behavior & Organization* 213:394–408.
- McCright, Aaron M., and Riley E. Dunlap. 2011. "The Politicization of Climate Change and Polarization in the American Public's Views of Global Warming, 2001–2010." *Sociological Quarterly* 52 (2): 155–194.
- Rabin, Matthew. 1995. *Moral Preferences, Moral Constraints, and Self-Serving Biases*. Technical report. Berkeley Department of Economics Working Paper No. 95-241.

Table C1: Overview over Statements and Charities.

	<b>Debate</b>	<b>Moral Statement</b>	<b>Factual Statement</b>	<b>Donation</b>
		<i>“How much do you agree with the following statement?”</i>	<i>“How likely do you think it is that the following statement is true?”</i>	Charity
1	Migration	People should be allowed to migrate freely between countries.	All countries benefit economically from the free movement of labour.	American Immigration Council
2	Animal Welfare	It is wrong to eat animals.	Animals feel less pain than humans.	World Animal Protection
3	Gender Equality	Gender equality should be an objective of policymaking.	Discrimination against women is the primary reason why women earn less than men.	Equality Now
4	Abortion	Abortion should be legal.	Women who have had an abortion experience more psychological distress than women who have had a miscarriage.	Planned Parenthood
5	Prostitution	Prostitution should be illegal.	Human trafficking is facilitated by liberal prostitution laws.	A21
6	Same-sex Marriage	Gay couples should have the same rights as heterosexual couples.	Societies where same-sex marriage is legal are happier than societies where it is illegal.	OutRight

Table C2: Description of Charities.

	<b>Debate</b>	<b>Charity</b>	<b>Text to introduce charity in experiment</b>
1	Migration	American Immigration Council	The American Immigration Council envisions an America that values fairness and justice for immigrants and believes that immigrants are part of the national fabric, bringing energy and skills that benefit all Americans. To advance change they engage in litigation, research, legislative and administrative advocacy, and communications.
2	Animal Welfare	World Animal Protection	World Animal Protection works towards a world where animals live free from suffering. They seek to improve the living conditions of animals farmed for food, to protect and save wild animals, animals affected by disasters, and working animals.
3	Gender Equality	Equality Now	Equality Now believes in creating a just world where women and girls have the same rights as men and boys. They use a unique combination of legal advocacy, regional partnership-building and community mobilization to encourage governments to adopt, improve and enforce laws that protect and promote the rights of women and girls.
4	Abortion	Planned Parenthood	The mission of Planned Parenthood is to provide comprehensive reproductive and complementary health care services in settings which preserve and protect the essential privacy and rights of each individual and to advocate public policies which ensure access to such services. They provide information and support to women considering to end a pregnancy in a clinic or using an abortion pill.
5	Prostitution	A21	The mission of A21 is to end human trafficking and slavery. They work closely with law enforcement on the ground to support police operations, identify victims through their hotlines, assist in the prosecution of traffickers, and represent survivors in court proceedings.
6	Same-sex marriage	OutRight	OutRight envisions a world where LGBTIQ (lesbian, gay, bisexual, transgender/transsexual, intersexual and queer) people everywhere enjoy full human rights and fundamental freedoms. They seek to fill research gaps, provide trainings to community members and allies to develop their expertise, and convene key stakeholders to exchange information on best practises related to ending violence based on sexual orientation.