

# Political Argumentation and Attitude Change in Online Interactions

Isaac D. Mehlhaff\*

August 29, 2022

## Abstract

Prevailing theories of public opinion and political psychology hold that human reasoning is biased and lazy, which suggests it is ill-suited to help ordinary citizens engage meaningfully with politics. In contrast, I draw on recent advances in cognitive psychology to contend that the biased and lazy nature of reasoning is precisely what gives citizens the tools to think through political issues and update their attitudes in response to argumentative exchanges. To test these hypotheses, I train a series of deep neural networks to classify textual inputs on several characteristics of discussion and argumentation. I use these classifiers to annotate over one million comments from the Reddit social media platform and show that attitude change is substantially more likely to result from argumentative exchanges rather than more contemplative ones. Results suggest that under the right conditions, humans can be quite skilled political reasoners.

---

\*The University of North Carolina at Chapel Hill; mehlhaff@live.unc.edu; word count: 8,987.

Citizens' ability to make reasoned choices about political issues and candidates is central to the function of democracy (Dahl 1998; Madison 1787; Madison 1788). This requisite of proper democratic practice has spawned extensive work on the question of whether citizens know enough to meaningfully participate in politics. Predominant theories of public opinion generally avoid attributing high-level reasoning skills to citizens (Farrell, Mercier, and Schwartzberg forthcoming), instead emphasizing attitude sources ranging from media and political elites (Achen and Bartels 2016; Zaller 1992) to heuristics such as party identification (Berelson, Lazarsfeld, and McPhee 1954; Lupia and McCubbins 1998). These perspectives suggest that citizens are unskilled reasoners; they are unable to logically think through political issues, so they rely on heuristics, stereotypes, and biases to make minimally informed decisions, often through a process of motivated reasoning.

I push back on this characterization. Drawing on recent advances in cognitive psychology, I argue that citizens' political reasoning capabilities are best utilized not when they attempt to reason on their own or even when they engage in casual political discussion. Instead, citizens are most likely to change their opinions and reduce their reliance on partisan stereotypes and heuristics when they engage in debate, exchanging and evaluating a series of arguments and counterarguments with an interlocutor. I evaluate this theory in the context of online political discussion. Approximately 70 percent of Americans use at least one social media platform (Pew Research Center 2021), making these types of online spaces one of the most common arenas in which citizens might encounter political content. I retrieve over one million comments posted on the Reddit social media platform over the last two years, and train several classifiers to annotate them for characteristics of argumentation. Results suggest consistent support for an argumentative theory of political reasoning: Attitude change is most likely to occur when, for example, interlocutors address each other directly and formulate high-quality counterarguments.

I make three main contributions in this paper. First, I develop a theory of the conditions under which we should expect reasoning to lead to attitude change on political issues. In doing so, I aim to turn the prevailing notion of political reasoning on its head and outline an approach to

understanding attitude change that takes the biased and lazy nature of human reasoning as an asset to leverage rather than an obstacle to overcome. Second, I demonstrate the inferential advantages of conducting research on Reddit, an underutilized data source in political science and psychology. Finally, I show how methodological advances in deep learning can assist scholars in understanding complex concepts in political speech.

## Political Reasoning and its Motivated Nature

For purposes of exposition, I contrast my theory of political reasoning below with the “contemplative” style of reasoning more common in political science and psychology. This idealized, *homo economicus* version of reasoning assumes that citizens can speak and think about political topics in a level-headed manner, are well-informed such that they can provide and assess evidence-based claims for or against a policy proposal, think and reason rationally, and are capable of arriving at the most logical and objectively optimal conclusion via this process of level-headed discussion and information consideration. It follows that the quality of conclusions should be monotonic increasing with respect to the amount and quality of information provided. Thus the “contemplative” label I apply to this framework: Citizens contemplate political information and use it to produce conclusions that are logically consistent with that information.

The empirical record, however, is mixed. Some scholars find that political discussions adhering to this contemplative model have positive impacts, increasing attitude strength, promoting consensus, and decreasing ideological polarization in the short-term (Becker, Porter, and Centola 2019; Druckman and Nelson 2003; Esterling, Fung, and Lee 2021). Others argue that they can actually exacerbate polarization or otherwise fail to promote more educated, well-reasoned opinions (Jackman and Sniderman 2006; Levendusky, Druckman, and McLain 2016; Mendelberg and Oleske 2000).<sup>1</sup> Indeed, this is the conclusion at which Taber and colleagues arrive following multiple studies. They argue that the processing of political arguments is motivated, and that both confirmation (seeking out confirmatory evidence) and disconfirmation (being skeptical of coun-

---

<sup>1</sup>On group polarization, see Arima (2012), Myers and Bishop (1970), and Sunstein (2002), among others.

terattitudinal arguments while uncritically accepting confirmatory arguments) biases both lead to attitude polarization (Taber, Cann, and Kucsova 2009; Taber and Lodge 2006). Results from studies investigating the effect of the media—the source of most political information in the present day—paint a picture of a citizenry that seeks out information reinforcing their existing beliefs if they are not pressured to do otherwise and rationalizes what little countervailing information they do consume to further support those beliefs (Hopkins 2014; Iyengar and Hahn 2009; Prior 2013). The potential for contemplative political reasoning to encourage attitude change or even nuanced consideration of diverse ideas appears weak.

Dual-process theories—the predominant paradigm for understanding the psychology of reasoning—suggest a deeper explanation for these findings. These theories separate cognition into Type 1 (quick, subconscious, and intuitive) and Type 2 (slow, conscious, and rational) and are often concerned with the ways in which preconscious processing can shape and distort conscious reasoning and decision-making (Evans and Stanovich 2013; Evans 2003; Kahneman 2011). Scholars find that the reasoning enabled by Type 2 cognition—supposedly conscious and rational—is instead prone to motivated reasoning (Ditto and Lopez 1992; Kunda 1990), frequently resorts to fast-and-frugal heuristics instead of slow-and-deliberate ratiocination (Chaiken 1980; Kahneman, Slovic, and Tversky 1982), and is barely more skilled at prediction and forecasting than the proverbial dart-throwing chimp (Mellers et al. 2015; Tetlock and Gardner 2015). These findings suggest that humans possess limited capacity for reasoning and the reasoning we do perform is unreliable.

The prognosis given by scholars studying political reasoning is similarly grim. Framing and priming effects abound (Bizer and Petty 2005; Chong and Druckman 2007; Tesler 2015), and individuals' political attitudes seem to be so sensitive to these effects that an influential strand of public opinion theory argues that most citizens lack coherent ideological principles in general (Kinder and Kalmoe 2017; Zaller and Feldman 1992; Zaller 1992). More perniciously, partisan and racial stereotypes are common (Ahler and Sood 2018; Gaertner and McLaughlin 1983), and they encourage citizens to reason in a politically motivated manner (Erisen, Lodge, and Taber 2014; Jost et al. 2003; Munro, Weih, and Tsai 2010). In a vein of the dual-process tradition, Lodge

and Taber (2000) develop a model of political reasoning that incorporates numerous heuristics and biases and asserts the central importance of hot cognition: the theory that reasoning is colored by an individual's affective state.<sup>2</sup>

In sum, many scholars argue that politics possesses limited immunity to the “flaws” in human reasoning posited in the psychology literature (Bolsen, Druckman, and Cook 2014; Gaines et al. 2007; Sniderman, Brody, and Tetlock 1991). Working in the predominant paradigm of political reasoning, then, one would conclude that citizens are poorly equipped to produce coherent reasoning on political issues and are resistant to changing or even forming sensible opinions.

## **An Argumentative Theory of Political Reasoning**

In contrast, I present a second framework, which I call “argumentative” political reasoning. This theory departs from existing theories in a major way: I consider it an empirical fact that human reasoning is biased and lazy. In the contemplative framework, these are *flaws* of reasoning because they complicate individuals' ability to cogitate rationally and objectively. In the argumentative framework, however, they are *features* of reasoning because they perform important functions.<sup>3</sup> This theory therefore takes as its point of departure the empirical literature that makes human reasoning look flawed and incorporates those findings as its core assumptions.

I define “argumentation” as giving, receiving, and responding to reasons in an iterative exchange. This type of interactive reason-giving likely approximates the environment in which human reasoning developed, and recent theoretical advances in cognitive psychology have pinpointed these interactions as critical for forming and updating attitudes (Mercier and Sperber 2011; Mercier and Sperber 2017). I contend that argumentation, an active manifestation of this ability to evaluate competing streams of information, is more effective than contemplation in changing political attitudes—and thereby ameliorating political polarization—because it pushes individuals to recognize the potential pitfalls of their own argument as well as the merits of different ideas. Evidence

---

<sup>2</sup>See also Taber, Cann, and Kucsova (2009) and Taber and Lodge (2016).

<sup>3</sup>See Haselton, Nettle, and Andrews (2005).

from psychology supports this assertion: Stronger, less fallacious arguments are more likely to lead to attitude change, even if they contradict prior beliefs (Hahn and Oaksford 2007; Koenig 2012; Trouche, Sander, and Mercier 2014).

The mechanisms linking argumentation to attitude change are the familiar forms of bias and laziness so thoroughly documented in the political psychology literature. Humans are cognitive misers by default; the brain is a hungry organ and will avoid expending energy unless pushed to do so (e.g. De Neys, Rossi, and Houdé 2013; Donald 1991). We are thus *lazy* in our generation of reasons to justify our opinions and behavior (Mercier and Landemore 2012). The optimal argument is the one that minimizes energy expenditure while maximizing the probability of our idea being accepted by others.<sup>4</sup> The ideal strategy to produce these arguments is therefore to offer the lowest-quality arguments first. If they are accepted, there is no need to spend time and energy coming up with more sophisticated ones. If they are rejected, we should gradually offer more and more complex arguments, addressing more and more of our interlocutor's rebuttals until one of those arguments is finally accepted. For similar reasons, we should expend relatively little energy on anticipating counterarguments, primarily because one argument can elicit a large number of potential counterarguments. Instead of anticipating and responding to all these counterarguments—especially when some might be unconvincing or things our interlocutor never would have mentioned in the first place—it is more efficient to instead use the iterative nature of argumentation to discern which types of arguments are effective, which are not, and which sorts of counterarguments must be addressed.

We are also *biased* in our evaluation of reasons given by others (Trouche, Johansson, et al. 2016). On average, ideas or actions supported by only low-quality arguments are less likely to be worth serious consideration compared to ones supported by high-quality arguments. In fact, *ceteris paribus*, the best idea is likely to be the one for which the best arguments can be generated. This implies that knowingly accepting anything less than the highest-quality arguments possible entails a higher likelihood of reaching a suboptimal conclusion. The optimal strategy, therefore, is

---

<sup>4</sup>On the tradeoff between information quality and energy expenditure, see Beach and Mitchell (1978), Payne, Bettman, and Johnson (1988), and Stigler (1961).

to only accept high-quality arguments that provide substantial reasons to believe in the veracity of our interlocutor's argument.

Evaluating high-quality reasons for an argument contrary to our predisposition while simultaneously being forced to generate high-quality reasons for our own argument are thus two components of the environment in which we should expect attitude change to occur.<sup>5</sup> Both dynamics must be present: Without someone to refute your own argument, you will not engage deeply with the reasons for holding your own opinions, and without challenging your interlocutor on their ideas, the reasons they give are unlikely to be convincing.<sup>6</sup> This is precisely the reason why engaging in contemplative political reasoning, consuming news media, or otherwise participating in intrapersonal consideration of political issues rarely promotes attitude change, even if multiple viewpoints are present.

Previous analyses of deliberative polling transcripts and laboratory experiments designed to evaluate these types of communication provide preliminary evidence for the persuasive power of argumentative reasoning compared to contemplative reasoning. Westwood (2015) is principally concerned with determining whether opinion change in group discussion of political issues is driven by reason-based interpersonal persuasion or information-based increases and refinements in political knowledge. Using transcripts from an online deliberative poll conducted in the United States,<sup>7</sup> he combines content and network analysis to model argument quality and the flow of those arguments between group members.

Consistent with an argumentative theory of political reasoning, results show that arguments presented to individuals in direct debate are of higher quality than those directed at the group more generally, high-quality arguments directed at individuals in direct debate are the strongest predic-

---

<sup>5</sup>One might be concerned that exposing individuals to counterattitudinal information would cause a backlash effect, leading them to adhere more tightly to their previous beliefs, but Guess and Coppock (2020) show across three survey experiments that this effect does not occur (see also Wood and Porter 2019).

<sup>6</sup>Another possible benefit to encouraging the generation of high-quality arguments is that it helps break the illusion of explanatory depth, wherein individuals believe they understand the nuances of an issue or policy but, when asked to explain it mechanistically, are forced to reckon with the fact that their knowledge is actually quite limited (Fernbach et al. 2013), cf. Crawford and Ruscio (2021).

<sup>7</sup>On deliberative polling, see Fishkin (2009), Luskin, Fishkin, and Jowell (2002), and Luskin, O'Flynn, et al. (2014).

tor of opinion change, and this opinion change occurs in the direction of the persuader's attitude (meaning that neither group polarization nor backlash effects occur in response to counterattitudinal arguments).<sup>8</sup> Moreover, political knowledge appears to have no bearing on opinion change, implying that high-quality arguments presented in direct debate with an interlocutor are persuasive regardless of either participant's level of political knowledge before or after the interaction.<sup>9</sup>

Schneiderhan and Khan (2008) contribute similar findings from a randomized, controlled laboratory experiment. Their first treatment arm, "discussion," entails free-flowing interpersonal conversation that typically involves the exchange of information but does not necessarily involve any form of reason-giving or argumentation. Their second treatment arm, "deliberation," is more akin to argumentation, where subjects are explicitly directed to process each other's opinions, provide their own opinions with justifications, and accept conflict in the discussion. Results show a clear advantage for the argument-style treatment. Subjects in this treatment are more likely to change their opinion compared to both the discussion treatment and the control, and the discussion treatment was statistically indistinguishable from the control on the matter of opinion change. The quality of argumentation also matters, as the number of justificatory reasons provided is positively associated with opinion change. Although the discussion topic (segregated student fees) was neither explicitly political nor highly contentious, these results provide additional evidence from a controlled setting for the role of argumentation in reasoning.

---

<sup>8</sup>Gerber et al. (2014) examine a separate deliberative poll conducted in Europe and arrive at similar, though slightly more ambivalent, conclusions.

<sup>9</sup>These findings comport with the theoretical argument made by Lupia (2002), who argues that the emphasis placed on information acquisition in deliberative environments is misguided. Left to their own devices, subjects presented with new—even balanced—information will merely rationalize that information in accordance with their predispositions.



## Data and Methods

### Reddit as a Research Tool: r/ChangeMyView

While experimental and survey methods are common currency in psychologically oriented studies, they leave questions as to their external validity (Barabas and Jerit 2010). These questions are even more critical in this case. Political discussion in the modern era is a social encounter (Carlson and Settle 2022); highly controlled experiments may remove many important elements of interpersonal interaction that we could capture by observing political argumentation “in the wild.”

Data from online media platforms is especially useful for gaining these insights (Baughan et al. 2021; Habernal and Gurevych 2017; Settle 2018); interactions are plentiful and data is frequently organized into comment-reply structures. I leverage Reddit data from a particular subreddit: r/ChangeMyView. This subreddit provides users a platform for engaging one another on questions both controversial and anodyne. An original poster (OP) begins a thread by stating their opinion on some topic, justifying that opinion, and inviting others to attempt to change all or part of that opinion. Other users then post comments and replies to either the original post or subsequent comments. This typically results in a deep comment forest with numerous users debating or agreeing with one another over specific points or broader arguments related to the topic.

Important for inferential purposes are the subreddit’s specific rules and thorough moderation practices. The subreddit is maintained by a dedicated group of moderators who remove posts and comments that do not adhere to the community’s rules. For an OP, these rules dictate that the initial explanation must thoroughly explain their opinion in no fewer than five hundred characters, they must genuinely hold that opinion (i.e. they are not playing devil’s advocate) and be open to it changing, their stance may not be neutral, and they must return to the post within three hours to extensively engage with their interlocutors. For non-OP users, comments given in direct reply to the original post must challenge the OP’s opinion with justification, they must refrain from making accusations of bad faith, and each comment must contribute meaningfully to the conversation (i.e.

they may not submit posts with only links, jokes, or non-substantive statements of agreement). All users are additionally held to common standards of online etiquette; they may not advocate for harm or make rude or hostile comments. Posts and comments not adhering to these rules are liberally removed by the moderators. The result of these moderation criteria is an open, interactive environment that ensures the final data collected from the platform bear close resemblance to the type of data one might expect to receive from a more controlled laboratory experiment, but without the constraints and inorganic interactions imposed by such a study.

The key benefit to this subreddit is its method of indicating opinion change. When a user—OP or otherwise—recognizes that part or all of their opinion has been changed by a particular post, they award that post a “delta.” The user awarding the delta must then explain which part of their opinion was changed by the post and why the post convinced them to change it. Moderators remove posts which award deltas sarcastically, as an expression of agreement, or for any purpose other than genuine opinion change. I use these deltas below to assess whether argumentation is more likely to result in attitude change relative to more common modes of discussion.

I draw these data from Pushshift’s Reddit data collection platform, an online infrastructure that automatically scrapes and archives all posts, comments, and associated metadata from Reddit. I pull two consecutive years of data from this archive, from July 1, 2020 to June 30, 2022. This time period encapsulates several contentious and politically relevant events from recent American history, with the 2020 presidential election cycle, the COVID-19 pandemic, and the Black Lives Matter protests perhaps the most significant. These data collection efforts yielded 12,593 posts, each with its own comment forest and discussion. The final comment-level dataset consists of 1,045,599 unique comments.

## **Detecting Characteristics of Argumentation**

To test the potential for argumentation to lead to attitude change in these online fora, I first need to determine whether each of these one million comments engages in debate or a more contemplative style of discussion. These data needs present a couple challenges. First, these types of

discussion characteristics can be complex, difficult to understand, and quite subjective. Identifying high-quality arguments generally involves asking survey respondents or experiment subjects whether they perceive an argument as strong or weak (Eagly and Chaiken 1993), resulting in “an inadequate tautology” that limits scholars’ ability to understand public opinion (Druckman 2022, p. 73). Instead of relying on one “argumentative or not” indicator—the coding of which could reflect multiple distinct ways of thinking about what argumentation looks like in practice—I come at this concept from several different angles, with each indicator capturing one narrow aspect of argumentative speech as summarized in Table 1. The result is a more complete picture of argumentation in online discourse.

Table 1: Summary of Argumentation Indicators

<b>Indicator</b>	<b>Description</b>	<b>Annotation Method</b>	<b>Training Data</b>
<b>Disagreement</b>	Does the commenter express disagreement or agreement?	Deep Neural Network	Internet Argument Corpus
<b>Object of Address</b>	Does the commenter direct their comment at a specific individual or the broader audience?	Quote Detection	None
<b>Question vs. Assert</b>	Does the commenter probe for more information or assert their own ideas?	Deep Neural Network	Internet Argument Corpus
<b>Counterargue vs. Rebut</b>	Does the commenter present an argument of their own or merely rebut a previous argument?	Deep Neural Network	Internet Argument Corpus
<b>Scope of Argument</b>	Does the commenter contradict the entirety of an interlocutor’s argument or just one specific part of the argument?	Deep Neural Network	Internet Argument Corpus
<b>Quality of Argument</b>	Does the comment express a clear, relevant, and high-impact argument?	Deep Neural Network	IBM-Rank-30k

Second, and more critically, training human annotators and paying them to read over one million documents is a process that would typically take months, cost thousands of dollars, and may not even result in reliable annotations despite those investments. Instead, I train a series of deep

neural networks to classify texts on five of the six characteristics summarized in Table 1. I conduct feature extraction for these classifiers with bidirectional encoder representations from transformers (BERT), a neural network architecture used in a wide variety of high-profile products such as Google Search (Devlin et al. 2019). BERT is pre-trained on English Wikipedia and the BooksCorpus (Zhu et al. 2015), which collectively provide a training corpus of over 3.3 billion words.

For four of these tasks, I draw training data from the Internet Argument Corpus, a collection of approximately 28,000 posts extracted from several online debate and discussion fora very similar to *r/ChangeMyView* (Abbott, Ecker, et al. 2016; Walker et al. 2012). The discussions in this corpus cover a variety of controversial topics relevant to politics and social life in the United States, such as same-sex marriage, gun control, and the existence of God. This diversity of issues is especially useful for training domain-general classifiers, as it prevents the models from over-fitting on words or phrases relevant to specific topics. For the fifth task, argument quality, I draw training data from IBM-Rank-30k, a corpus of approximately 30,000 crowd-sourced arguments across a similarly diverse set of 71 common topics (Gretz et al. 2020). Both sets of training data are coded on each characteristic by 5-10 human annotators. Additional details on training data are provided in the Supplementary Information.

Figure 1 displays accuracy and weighted F1 scores for each of the five classifiers trained to annotate Reddit comments on these characteristics, comparing the performance of the deep neural networks to two baselines. The naïve baseline represents the performance of a classifier that merely selects a class at random, while the lexical baseline contrasts the deep learning approach with a bag-of-words approach that is more common in political science. Focusing in particular on the weighted F1 scores—which are preferred to accuracy metrics when classes are imbalanced, as they are here—reveals that classifying characteristics of argumentation is a challenging task. Nevertheless, the deep neural networks perform well, consistently besting the lexical baseline. F1 scores above 0.7 are strong for this type of task, and the deep neural network classifier even achieves state-of-the-art results on the disagreement identification task (Abbott, Walker, et al. 2011; Wang and

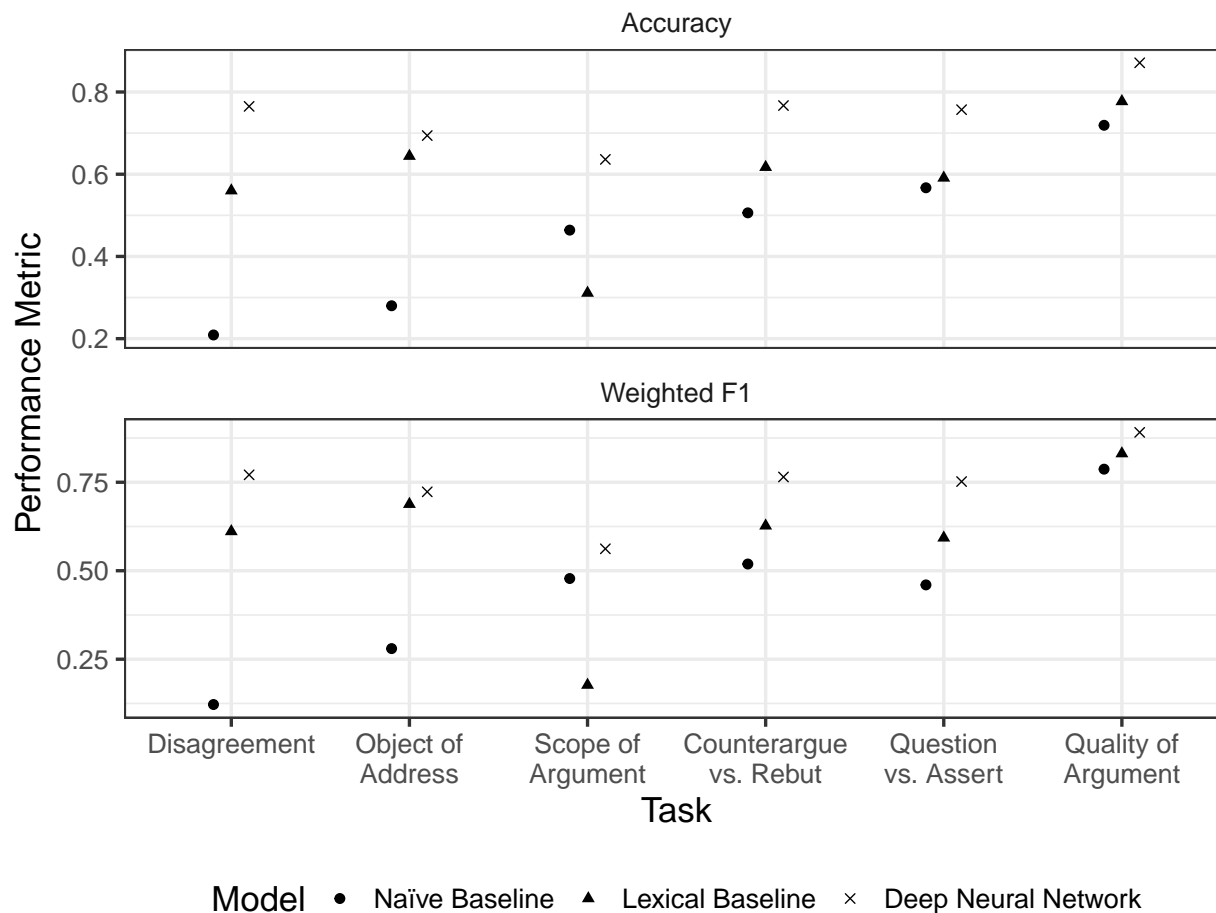


Figure 1: Accuracy and Weighted F1 Scores. Naïve baselines randomly select classes, lexical baselines use support vector machines with stochastic gradient descent and features extracted from word frequency arrays. See Supplementary Information for additional details on model training and testing.

Cardie 2014).<sup>10</sup> These results suggest strong potential for deep learning methods to assist scholars in understanding complex phenomena—such as argumentation—in political speech.

Finally, one indicator—object of address—does not draw its annotations from a deep learning model. This indicator captures whether a commenter is responding to or engaging with an interlocutor directly, as opposed to making comments for a broader audience. The structure of the r/ChangeMyView text data allows me to detect this directly. Commenters can respond directly to previous comments by “quoting” them. These quotes are copied into the commenter’s new post

<sup>10</sup>Additional performance metrics, details on model training, and tests of five other classifier architectures are provided in the Supplementary Information.

in a text box set off from the rest of the post. If a commenter quotes a previous post, I code that comment as addressing an individual. If not, I code it as being meant for the audience as a whole.

## Drawing Inferences from Machine-Proxied Data

I noted above the benefits conferred by using machine learning methods to annotate text data. However, these methods also impose an additional challenge. The annotations they produce are only *proxies*, denoted here by  $\hat{X}$ ; the true values of these indicators  $X$  are unobserved and, in many cases, unobservable. Using  $\hat{X}$  to proxy for  $X$  in inferential statistical analyses requires special care, as it causes two problems.

First, using learned proxies as explanatory variables leads to attenuation bias, pushing coefficient estimates toward zero (Wooldridge 2015). However, unbiased estimates of the effect of  $X$  on  $Y$  can still be drawn under three assumptions (VanderWeele and Hernán 2012):<sup>11</sup>

1. **Positive average monotonicity:** When  $X$  is larger,  $\hat{X}$  is also larger, on average. This is an especially weak assumption when using supervised learning methods, as the relationship between  $\hat{X}$  and  $X$  can be empirically assessed. Performance metrics in Figure 1 and the Supplementary Information suggest that this assumption holds.
2. **Perfect observation of outcome and covariates:** The outcome  $Y$  and covariates  $W$  are perfectly observed. This is likely to hold in this case, as  $Y$  and  $W$  are not proxied and instead represent directly observable quantities available in Reddit metadata.
3. **Specification:** The covariates  $W$  are correctly specified. This is the strongest assumption in this context, but the anonymity of Reddit ameliorates many specification concerns; although it is impossible to adjust for variables that often affect political discussion—such as race, gender, or partisanship—discussion participants also lack this information about each other, decreasing the likelihood that they will exert an effect on either  $\hat{X}$ ,  $X$ , or  $Y$ .

<sup>11</sup>Even under these assumptions, it is not necessarily true that a failure to find an effect suggests the lack of such an effect, for the familiar reasons of statistical power and the nature of null-hypothesis significance testing (Knox, Lucas, and Cho 2022). With such a large dataset, however, concerns about power are likely unfounded in this case.

Second, using learned proxies will result in downward-biased standard errors, as results from the generalized linear models that I present below incorporate two sources of sampling variability. The first source of sampling variability comes from drawing one of many possible datasets for training the deep neural networks to estimate  $\hat{X}$ . The second comes from drawing a sample of observations from a larger population of possible observations to estimate the effect of  $\hat{X}$  on  $Y$ . That is, uncertainty enters the empirical analysis in both the measurement model and the linear model, and both sources must be accounted for.

To do so, I follow Knox, Lucas, and Cho (2022), who recommend a bootstrap approach to estimating coefficient estimates with appropriate standard errors when the key explanatory variable is a learned proxy.<sup>12</sup> This bootstrap has two stages. In the first stage, I randomly resample the training set with replacement 500 times, re-train the deep neural network with each resampled training set, and estimate  $\hat{X}$ . In the second step, I randomly resample the final Reddit dataset with replacement 500 times, re-fitting the substantive models below with each resampled dataset. This procedure results in 250,000 total coefficient estimates, from which accurate standard errors can be calculated.<sup>13</sup>

## Drawing Inferences from Rare-Events Data

One final methodological consideration is in order before moving to substantive results. The dependent variable, attitude change, is represented here as a binary phenomenon—either a comment succeeds in changing someone’s opinion or it does not. I thus rely on binomial logits to estimate the effect of argumentative characteristics on attitude change. Attitude change, however, is a rare event; descriptive analyses presented below show that only 1.7 percent of all r/ChangeMyView comments are awarded a delta, and 36.2 percent of all comment forests award no deltas at all. Lo-

<sup>12</sup>Though not directly applicable, see also Fong and Grimmer (forthcoming) on estimating causal effects with latent treatments.

<sup>13</sup>Having well over one million observations in a linear model may lead to concerns that any statistically significant findings are merely the result of an extremely large sample size, so I sample only 10 percent of the total number of observations in the second step prior to fitting each iteration of the models. The Supplementary Information displays results for the full dataset as well as results from post-level sampling instead of comment-level sampling.

gistic regression on such unbalanced data leads to downward-biased estimates of event probabilities (King and Zeng 2001). I therefore adopt a penalized maximum-likelihood approach to correct for this bias. Firth (1993) proposed using the Jeffreys prior—the square root of the determinant of the Fisher information matrix—to penalize the log-likelihood function. This penalty has been shown to reduce bias in logistic coefficient estimates in the cases of small samples, rare events, and complete or quasi-complete separation (Heinze and Schemper 2002; Puhr et al. 2017). All results derived from binomial logits below use this penalized maximum-likelihood framework.

## How Do People Discuss Politics Online?

This section presents a brief descriptive analysis of the `r/ChangeMyView` data and what the deep learning classifiers reveal about how citizens discuss contentious issues in online spaces. I turn first to summary data on how discussions unfold. Figure 2 contains four plots with information on post-level (i.e. discussion-level) characteristics. Plot A displays the number of deltas awarded in each post, plot B displays the number of unique individuals participating in the discussion, and plots C and D display the number of comments overall as well as the number of comments made by the OP. Red vertical lines locate the mean of each distribution, and labels on each plot report the mean, median, standard deviation, and range.

These aggregate statistics suggest a couple key takeaways. First, online discussions in this context appear similar to political discussions that occur face-to-face (see, for example, Carlson and Settle 2022). Distributions in Figure 2, plots C and D are heavily right-skewed, reflecting the fact that most discussions are relatively brief. The median post contains 51 comments—12 of which are by the OP—and these comments are typically spread among several unique threads. Figure 2, plot B is also heavily right-skewed, suggesting that these discussions typically occur in small- to medium-sized groups, with a median of 20 people participating. Typically, each person in these groups contributes only a few points to the discussion before leaving, but there is a small cadre of enthusiastic, heavily involved commenters. The mean number of comments contributed to



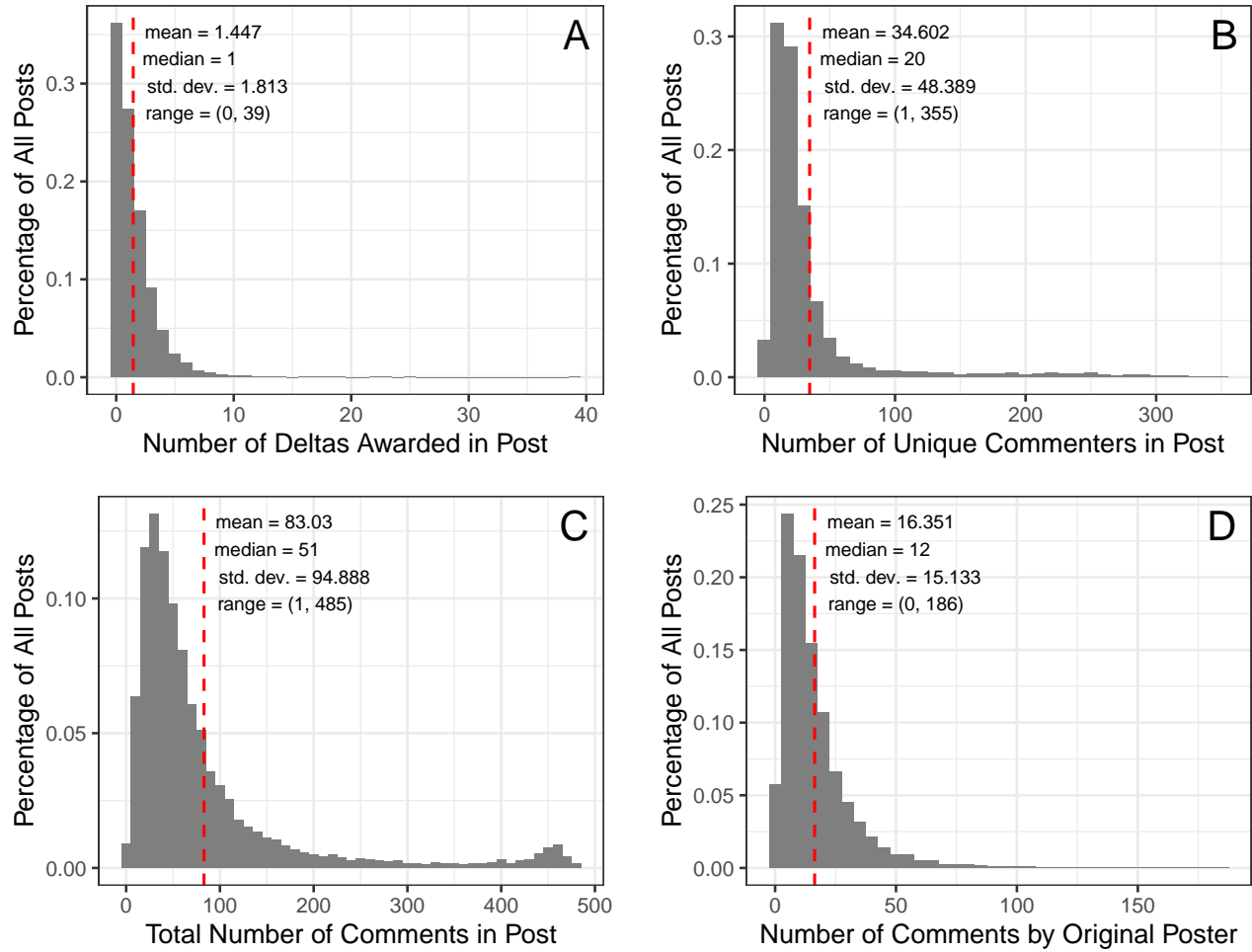


Figure 2: Descriptive Statistics on Discussion Participants and Deltas. Vertical red line shows mean of each distribution.

one post by a commenter is approximately 2.4, but the observed range runs as high as 186. These broad similarities to face-to-face political discussion imply that the substantive findings below should carry some external validity.

Second, r/ChangeMyView rules stipulate that OPs must be open to their opinion changing, but Figure 2, plot A suggests that attitude change is still not a common occurrence. In fact, the modal outcome is for no deltas to be awarded anywhere in a comment forest. Narrowing the focus to only those deltas awarded by OPs, the likelihood of a delta being awarded drops further: 38.8 percent of posts never receive a delta from the OP, suggesting that nobody was able to successfully change even part of the OP's opinion. Assessed as a response to individual arguments, opinion

change appears even rarer: Only 1.7 percent of comments are awarded a delta, and only 1.6 percent of comments are awarded a delta by an OP. In addition to ensuring adequate variation on this study's dependent variable, this finding alleviates concerns about selection bias. Although some participants may have come to these discussions more open-minded than the average politically engaged citizen, it does not appear that they were especially likely to change their opinion.

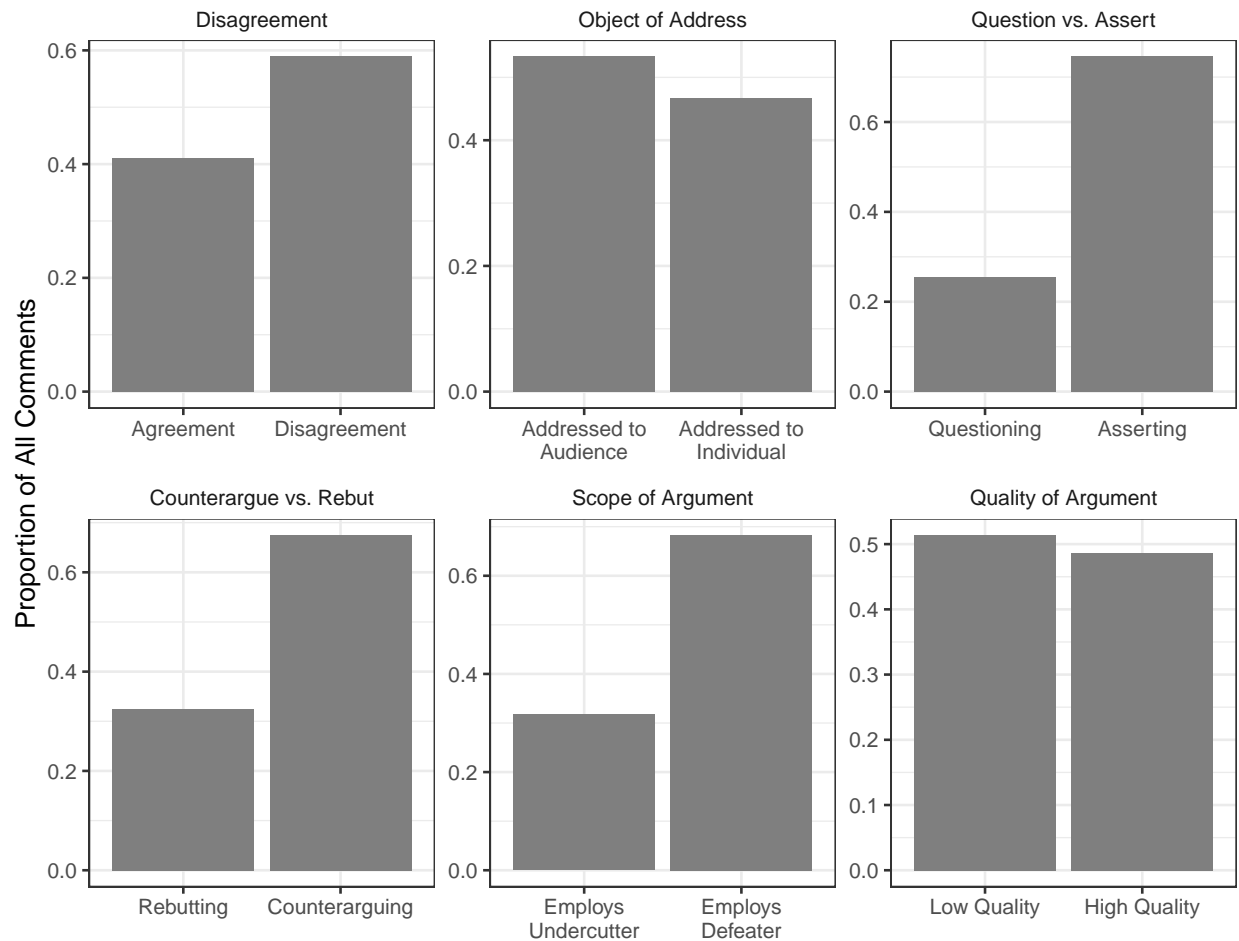


Figure 3: Class Frequencies of Argumentation Characteristics. Frequencies broken down by commenter type are presented in the Supplementary Information.

Figure 3 shows class frequencies for each of my key explanatory variables. Each exhibits substantial variation, providing reassurance that r/ChangeMyView is not merely hosting debates where every comment is a counterargument. Indeed, over 40 percent of comments are coded as expressing agreement with a previous comment, and approximately 25 percent of comments

merely probe for more information rather than making any assertions of their own. In addition, slightly over half of all comments are *not* directed at a specific interlocutor, indicating that the majority of these interactions are not characterized by back-and-forth argumentation.

## Can Argumentation Lead to Attitude Change?

I now turn to explicitly testing key theoretical implications. I focus on the six characteristics of argumentation described in Table 1 and displayed in Figure 3: disagreement, object of address, question or assertion, rebuttal or counterargument, and the scope and quality of arguments. Each of these variables taps a slightly different concept related to argumentation and, taken together, they paint a comprehensive portrait of the effect of argumentation on attitude change.

All models presented below use penalized maximum-likelihood to fit binomial logits, and all include three covariates. First, next to their username, *r/ChangeMyView* displays the number of deltas each commenter has been awarded over the course of their entire tenure on the subreddit. I control for this value to account for the possibility that discussion participants may view this record of deltas as a source of credibility and therefore be more likely to be swayed by commenters who have been awarded many deltas in the past. Second, the location of a comment in the comment forest may affect the likelihood that it receives a delta. In particular, the deeper in the comment forest a comment appears, the less likely that discussion participants will see that comment and award it a delta. I therefore control for the depth of each comment in the comment forest. Third, like most other subreddits, *r/ChangeMyView* allows users to up-vote or down-vote comments, and the balance of these votes is displayed next to each comment. While this voting process does not indicate attitude change, discussion participants may take the relative balance of those votes as evidence of a comment's quality (or lack thereof) and be more (or less) likely to be swayed as a result, so I include the comment's overall score as a control.

## Effects of Argumentation Characteristics on Probability of Attitude Change

Figure 4 displays the predicted probability of a comment resulting in attitude change, conditional on each characteristic of argumentation, which are each presented in a separate facet.<sup>14</sup> To place results in context, horizontal red lines on each facet show the baseline probability that any given comment will result in attitude change, conditional on the same covariates. Error bars denote 95 percent confidence intervals.

I begin with the broad hypothesis that higher levels of disagreement should be more likely to lead to attitude change. Disagreement itself does not imply argumentation, but it is likely a necessary element of argumentation; it makes little sense to critique a series of counterarguments if the interlocutors hold similar attitudes. Indeed, the idea that individuals are likely to uncritically accept statements with which they agree is one of the central findings of research on motivated reasoning (e.g. Bolsen, Druckman, and Cook 2014; Lebo and Cassino 2007; Stanley, Henne, et al. 2020). It is also an important assumption for the theory I present above. This test therefore serves as both a first cut at gauging the feasibility of the theory, as well as a validation check of sorts. Consistent with previous literature, comments expressing agreement are less likely to result in attitude change compared to those expressing disagreement. Although the difference in predicted probabilities is not statistically significant at the  $p < 0.05$  level, the coefficient estimate on the disagreement indicator (presented in Figure 5 below) is statistically different from zero, suggesting that individuals are more likely to change their opinions when exposed to counterattitudinal views.

Recall that a key definitional component of argumentation is an iterative exchange where interlocutors can process and respond directly to each other's arguments. Comments in which an individual responds directly to another comment and addresses their counterargument toward an individual should therefore be more likely to result in attitude change than comments that are meant for a general audience, even if those comments present high-quality arguments. The second (top middle) facet in Figure 4 suggests this is the case; comments directed at a specific discussion partner are significantly more likely to lead to attitude change, increasing the predicted probab-

---

<sup>14</sup>Note that each characteristic is evaluated in a separate model.

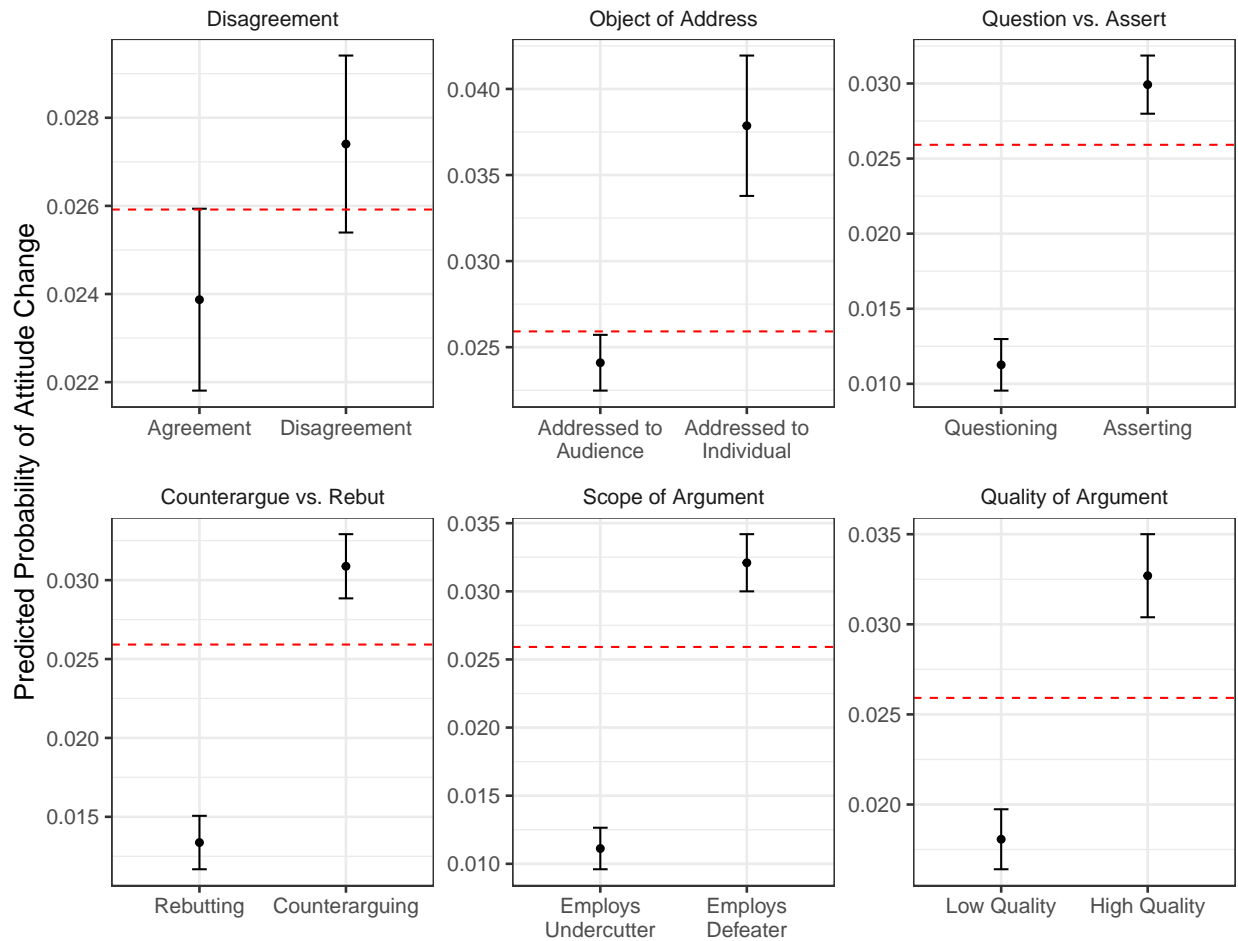


Figure 4: Effect of Argumentation Characteristics on Attitude Change. Red horizontal lines denote baseline probability of a comment resulting in attitude change. “Participants” make at least two comments. “Lurkers” only comment to award a delta. Error bars give 95% confidence intervals. Full results available in the Supplementary Information.

ity from 2.4 to 3.8 percent. While the magnitude of these probabilities is not substantial, it bears reiterating that attitude change is a rare event, with just 1.7 percent of all comments successfully achieving persuasion. In that sense, a relative increase of 58.3 percent is an appreciable effect size.

The next three tests are concerned with how the *type* of comment affects its persuasiveness. Theoretically, I expect comments to be more likely to lead to attitude change if they assert new ideas instead of probing for information, present original counterarguments instead of merely rebutting other arguments, and use counterarguments that address the entirety of an argument (i.e. a defeater) instead of only critiquing one piece of the argument (i.e. an undercutter). Results of

these indicators are presented in the top right, bottom left, and bottom middle facets, respectively, of Figure 4. Not only do results conform to expectations, but it is on these indicators that I observe the strongest effects.

Asserting new ideas is associated with a near-tripling of the predicted probability of attitude change compared to questioning an interlocutor, suggesting that the value of the Socratic method may not extend to persuasive appeals (Meckstroth 2012). This indicator provides perhaps the closest fit with the theoretical distinction between contemplative and argumentative reasoning, and results suggest that argumentative styles of discussion are substantially more likely to result in attitude change. Similarly, presenting an actual counterargument more than doubles the likelihood that a comment will lead to attitude change compared to merely deflecting other arguments. This is another critical test, as the counterargue/rebut indicator captures a nuanced distinction between argument types. That is, rebuttals might themselves be considered a type of argument, but they do not provide an affirmative reason to adopt a different viewpoint, only a negative reason to reject the opinion their interlocutor already holds. As I detailed in the theory section above, this approach is not likely to be successful, and results corroborate this expectation. Finally, the scope of arguments appears to have a dramatic effect on the likelihood of a comment resulting in attitude change. Arguments aimed at disproving an entire belief (defeaters) are approximately three times as likely to result in attitude change than arguments aimed at removing evidentiary support for a belief (undercutters).

Finally, I examine the effect of argument quality in the last (bottom right) facet of Figure 4. If argumentation had no effect on attitudes, I would expect to see no difference in the persuasiveness of low- or high-quality arguments; the quality of arguments would be irrelevant if argumentation itself was ineffective. Instead, putting forth a high-quality argument is associated with a relative increase of 83.3 percent in the probability of inducing attitude change—moving from 1.5 to 3.3 percent compared to low-quality arguments. Taken together, these results provide support for the argumentative theory of political reasoning across a wide range of indicators, with both statistically and substantively significant effects.

## Level of Engagement as a Moderator

Most people who view online fora are “lurkers;” they observe what others are discussing but do not (or only very rarely) contribute to the conversation themselves. Some estimates suggest the proportion of users falling into this category may be as high as 90 percent (Lukin et al. 2017). This feature of online discussion offers an inferential advantage in the context of the argumentative theory of reasoning. OPs—those who engage most in the discussion—should be more affected by the degree of argumentation than are lurkers, who may be less likely to experience the activation of cognitive processes enabling attitude change due to their lack of engagement. To evaluate this possibility, I distinguish between deltas awarded by the post’s OP, participants (users who comment at least one other time in addition to awarding a delta), and lurkers (users whose only contribution to the discussion is the comment in which they award a delta).<sup>15</sup>

To gauge the relative size and significance of effects, Figure 5 presents the estimated coefficients on each argumentation indicator, separated by the type of user awarding the delta. Error bars show 95 percent confidence intervals, and the red horizontal line denotes zero. As expected, effects are generally stronger among OPs compared to participants and lurkers. Coefficient estimates for the former are always positive and statistically significant, while those for the latter are sometimes negative, only statistically significant in two cases, and always carry point estimates lower than those for OPs. I take these consistent results as evidence of a moderation effect. One mechanism connecting argumentative interactions to attitude change is the tailoring of counterarguments in response to an interlocutor’s arguments, and the close examination of one’s own beliefs that those counterarguments prompt. Individuals not participating fully in argumentation are therefore less likely to reap the benefits of them, as they are not prompted to closely engage with the reasons why they hold their own opinions, nor is the discussion likely to generate arguments that speak directly to those reasons.

---

<sup>15</sup>The Supplementary Information presents class frequencies like those presented in Figure 3 broken down by commenter type.

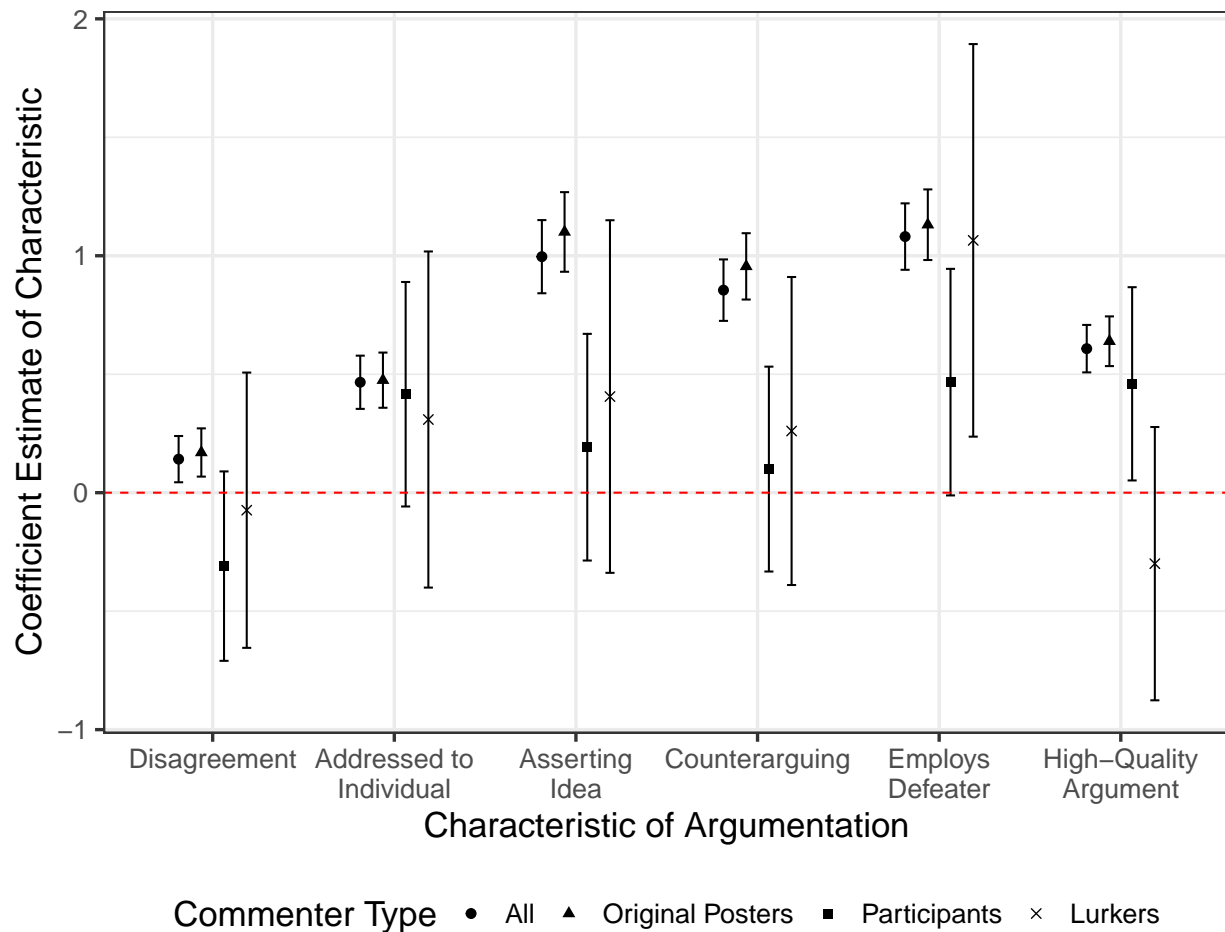


Figure 5: Coefficient Estimates of Argumentation Characteristics. Red horizontal line denotes zero. “Participants” make at least two comments. “Lurkers” only comment to award a delta. Error bars give 95% confidence intervals. Full results available in the Supplementary Information.

## Discussion

I have argued that the function of reasoning is to generate reasons justifying one’s behavior, evaluate reasons given by others, and determine when and to what extent attitude change is warranted based on the reasons given. The environment in which we should expect attitude change to occur, then, is one that encourages the production of high-quality reasons justifying one’s attitude and high-quality reasons rebutting that attitude. Finally, political argumentation is perhaps the only form of interpersonal interaction that naturally constructs an environment with both of these components.



Consider how this form of reasoning contrasts with the contemplative style of reasoning more common throughout the political psychology literature, and just how narrow is the range of interactions encompassed by argumentative reasoning. At minimum, argumentation requires two individuals exchanging arguments and counterarguments in iterative fashion, directly responding to one another's reasons. Simply watching a debate—even if the debaters meet this requirement—or reading sets of pro and con arguments still fall into the “contemplative” category because they do not invoke direct, two-sided cognitive engagement.

Given the biased and lazy nature of human reasoning, it is not surprising that these forms of information consumption can even lead to attitude polarization, as they do in Taber, Cann, and Kucsova (2009) and Taber and Lodge (2006).<sup>16</sup> In these studies, subjects do not have their beliefs challenged, their own attitudes are partially validated by being presented with reasons favoring those attitudes, and there is virtually no barrier to rationalizing away the counterattitudinal information. Presenting subjects with only counterattitudinal arguments, however, can elicit some attitude change. Gibson (1998) and Sniderman and Piazza (1993) attempt to persuade survey respondents by rebutting their views on political and racial tolerance, respectively. After presenting counterattitudinal arguments, Gibson observes up to 23.5 percent of respondents becoming more politically tolerant and Sniderman and Piazza observe up to 44 percent of respondents becoming more supportive of racial policies—strong results suggesting that presenting counterattitudinal arguments in direct response to a stated attitude can indeed lead individuals to reconsider that attitude, even on hotly contested, moralized issues. The results I present above add to this empirical pattern.

Another reason that non-argumentative studies may lead to attitude polarization is that the quality of arguments in these settings is likely quite low, because presenting subjects with a pre-determined set of reasons makes it impossible to address counterarguments in a flexible manner. When presented with a low-quality argument, it is easier for subjects to find fatal problems in the argument such that their confidence in the veracity of their own attitude is actually bolstered.

---

<sup>16</sup>See also Ansolabehere and Iyengar (1995), Redlawsk (2002), and Stanley, Henne, et al. (2020).

This effect is even more pronounced in individuals of high political sophistication, as they have greater domain-specific knowledge and are able to generate more counterarguments. This is why presenting high-quality counterarguments is so critical. It is also why argumentation is necessary; high-quality arguments calibrated to meet the concerns of each individual are unlikely to be produced in any other environment precisely due to the lazy nature of reason-giving.

I suggest that political argumentation provides at least three additional benefits that future research should assess. The first of these benefits is cognitive and the other two stem from the social nature of political argumentation. First, conditional on the provision of high-quality arguments, participants should become less confident in the veracity of their opinion, they should hold that opinion less strongly, and they should be more willing to compromise in the future, even if their actual preferences or ideal points do not move. Second, by pushing participants to search for better and more refined reasons to support their predispositions, it should lead to increased epistemic quality of arguments.<sup>17</sup> Finally, it exposes participants to real out-group partisans presenting real out-group partisan arguments; encroaching on partisan echo chambers and isolating out-group partisan arguments from media influences may humanize out-group partisans and decrease the existential threat they are considered to pose to one's in-group. The opposing viewpoint may seem accessible rather than foreign and each side may feel like their voices were heard.<sup>18</sup>

This final prediction may seem a lofty goal in a polarized political environment, but two recent studies provide empirical support. Stanley, Whitehead, et al. (2020) found that partisans believed their counter-partisans were less likely to have good reasons for their political attitudes, and that this expected lack of high-quality reasons spilled over into doubts about counter-partisans' intellectual and moral fortitude. Exposing these partisans to counterattitudinal arguments, however, led them to produce more favorable views of the counter-partisans who produce those arguments. Moreover, this effect was entirely independent of persuasion, suggesting that even if individuals do not change their attitudes as a result of argumentation, simply being exposed to high-quality ar-

---

<sup>17</sup>A significant amount of work on epistemic quality has been conducted in the literature on science education (e.g. Erduran and Jiménez-Aleixandre 2007; Kuhn 1992; Sandoval 2003).

<sup>18</sup>This may partially depend on interlocutors providing each other with "high-quality listening" (Itzhakov, Kluger, and Castro 2017; Kalla and Broockman 2020) in addition to high-quality arguments.

guments by their interlocutors may decrease the negative affect they feel toward those with whom they disagree.<sup>19</sup>

Dorison, Minson, and Rogers (2019) explicate the link between information search and partisan affect, showing that partisans' unwillingness to voluntarily expose themselves to counterattitudinal information is partially due to overestimating the strength of negative affect they are likely to feel in response to that information. More importantly, correcting that "affective forecasting error" led to greater voluntary exposure to counterattitudinal arguments. Combining these findings with those of Stanley, Whitehead, et al. (2020) suggest a positive feedback loop: Exposure to counterattitudinal arguments via argumentation leads to a decrease in negative partisan affect,<sup>20</sup> this decrease in negative partisan affect leads to increased willingness to engage with counterattitudinal arguments, and so on. Parsons (2010) uses observational data to test these effects and comes away with a clear result: Exposure to political disagreement depolarizes party affect.

In sum, argumentation appears to hold promise for changing political attitudes. Online discussions adhering to an array of argumentative characteristics are substantially more likely to result in attitude change compared to discussions with more contemplative characteristics. These findings comport with a wide range of related, though theoretically and methodologically distinct, studies showing the value of an interactionist approach to cognitive psychology and an argumentative theory of reasoning. In my view, this suggests reason for optimism. Under the correct conditions, humans can, in fact, be quite skilled political reasoners.

## References

Abbott, Rob, Brian Ecker, et al. (May 2016). "Internet Argument Corpus 2.0: An SQL Schema for Dialogic Social Media and the Corpora to Go With It". In: *Proceedings of the Tenth In-*

---

<sup>19</sup>It may also be the case that negative partisan affect will decrease as a downstream effect of decreased ideological distance, but Stanley, Whitehead, et al. (2020) do not test this directly.

<sup>20</sup>It should be emphasized that negative partisan affect is distinct from anxiety or negative core affect more generally. The former describes conscious dislike of opposing partisans while the latter two describe internal emotional and affective states. Negative partisan affect should therefore not be construed to regulate information search, openness to persuasion, or persuasiveness as argued above for anxiety and negative affect.

- ternational Conference on Language Resources and Evaluation*. Portorož, Slovenia, pp. 4445–4452.
- Abbott, Rob, Marilyn Walker, et al. (June 2011). “How Can You Say Such Things?!? Recognizing Disagreement in Informal Political Argument”. In: *Proceedings of the Workshop on Language in Social Media*. Portland, OR: Association for Computational Linguistics, pp. 2–11.
- Achen, Christopher H. and Larry M. Bartels (2016). *Democracy for Realists: Why Elections Do Not Produce Responsive Government*. Princeton, NJ: Princeton University Press.
- Ahler, Douglas J. and Gaurav Sood (July 2018). “The Parties in Our Heads: Misperceptions about Party Composition and Their Consequences”. In: *The Journal of Politics* 80.3, pp. 964–981.
- Ansolabehere, Stephen and Shanto Iyengar (1995). *Going Negative: How Political Advertisements Shrink and Polarize the Electorate*. New York: The Free Press.
- Arima, Yoshiko (Apr. 2012). “Effect of Group Means on the Probability of Consensus”. In: *Psychological Reports* 110.2, pp. 607–623.
- Barabas, Jason and Jennifer Jerit (May 2010). “Are Survey Experiments Externally Valid?” In: *American Political Science Review* 104.2, pp. 226–242.
- Baughan, Amanda et al. (Apr. 2021). “Someone Is Wrong on the Internet: Having Hard Conversations in Online Spaces”. In: *Proceedings of the ACM on Human-Computer Interaction* 5 (CSCW1), pp. 1–22.
- Beach, Lee Roy and Terence R. Mitchell (July 1978). “A Contingency Model for the Selection of Decision Strategies”. In: *The Academy of Management Review* 3.3, pp. 439–449.
- Becker, Joshua, Ethan Porter, and Damon Centola (May 2019). “The Wisdom of Partisan Crowds”. In: *Proceedings of the National Academy of Sciences* 116.22, pp. 10717–10722.
- Berelson, Bernard R., Paul F. Lazarsfeld, and William N. McPhee (1954). *Voting: A Study of Opinion Formation in a Presidential Campaign*. Chicago: The University of Chicago Press.
- Bizer, George Y. and Richard E. Petty (Aug. 2005). “How We Conceptualize Our Attitudes Matters: The Effects of Valence Framing on the Resistance of Political Attitudes”. In: *Political Psychology* 26.4, pp. 553–568.

- Bolsen, Toby, James N. Druckman, and Fay Lomax Cook (June 2014). “The Influence of Partisan Motivated Reasoning on Public Opinion”. In: *Political Behavior* 36.2, pp. 235–262.
- Carlson, Taylor N. and Jaime E. Settle (2022). *What Goes Without Saying: Navigating Political Discussion in America*. New York: Cambridge University Press.
- Chaiken, Shelly (Nov. 1980). “Heuristic Versus Systematic Information Processing and the Use of Source Versus Message Cues in Persuasion”. In: *Journal of Personality and Social Psychology* 39.5, pp. 752–766.
- Chong, Dennis and James N. Druckman (Nov. 2007). “Framing Public Opinion in Competitive Democracies”. In: *American Political Science Review* 101.4, pp. 637–655.
- Crawford, Jarret T. and John Ruscio (Apr. 2021). “Asking People to Explain Complex Policies Does Not Increase Political Moderation: Three Preregistered Failures to Closely Replicate Fernbach, Rogers, Fox, and Sloman’s (2013) Findings”. In: *Psychological Science* 32.4, pp. 611–621.
- Dahl, Robert A. (1998). *On Democracy*. 2nd ed. New Haven, CT: Yale University Press.
- De Neys, Wim, Sandrine Rossi, and Olivier Houdé (Apr. 2013). “Bats, Balls, and Substitution Sensitivity: Cognitive Misers Are No Happy Fools”. In: *Psychonomic Bulletin & Review* 20.2, pp. 269–273.
- Devlin, Jacob et al. (May 2019). “BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding”. Pre-Print. Google AI Language. arXiv: 1810.04805.
- Ditto, Peter H. and David F. Lopez (Oct. 1992). “Motivated Skepticism: Use of Differential Decision Criteria for Preferred and Nonpreferred Conclusions”. In: *Journal of Personality and Social Psychology* 63.4, pp. 568–584.
- Donald, Merlin (1991). *Origins of the Modern Mind: Three Stages in the Evolution of Culture and Cognition*. Cambridge, MA: Harvard University Press.
- Dorison, Charles A., Julia A. Minson, and Todd Rogers (July 2019). “Selective Exposure Partly Relies on Faulty Affective Forecasts”. In: *Cognition* 188, pp. 98–107.

- Druckman, James N. (2022). “A Framework for the Study of Persuasion”. In: *Annual Review of Political Science* 25.
- Druckman, James N. and Kjersten R. Nelson (Oct. 2003). “Framing and Deliberation: How Citizens’ Conversations Limit Elite Influence”. In: *American Journal of Political Science* 37.3, pp. 729–745.
- Eagly, Alice H. and Shelly Chaiken (1993). *The Psychology of Attitudes*. Fort Worth, TX: Harcourt Brace Jovanovich.
- Erduran, Sibel and María Pilar Jiménez-Aleixandre, eds. (2007). *Argumentation in Science Education: Perspectives from Classroom-Based Research*. Science & Technology Education Library 35. New York: Springer.
- Erisen, Cengiz, Milton Lodge, and Charles S. Taber (Apr. 2014). “Affective Contagion in Effortful Political Thinking”. In: *Political Psychology* 35.2, pp. 187–206.
- Esterling, Kevin M., Archon Fung, and Taeku Lee (Apr. 2021). “When Deliberation Produces Persuasion Rather than Polarization: Measuring and Modeling Small Group Dynamics in a Field Experiment”. In: *British Journal of Political Science* 51.2, pp. 666–684.
- Evans, Jonathan St. B. T. and Keith E. Stanovich (May 2013). “Dual-Process Theories of Higher Cognition: Advancing the Debate”. In: *Perspectives on Psychological Science* 8.3, pp. 223–241.
- Evans, Jonathan St. B.T. (Oct. 2003). “In Two Minds: Dual-Process Accounts of Reasoning”. In: *Trends in Cognitive Sciences* 7.10, pp. 454–459.
- Farrell, Henry, Hugo Mercier, and Melissa Schwartzberg (forthcoming). “Analytical Democratic Theory: A Microfoundational Approach”. In: *American Political Science Review*.
- Fernbach, Philip M. et al. (June 2013). “Political Extremism Is Supported by an Illusion of Understanding”. In: *Psychological Science* 24.6, pp. 939–946.
- Firth, David (Mar. 1993). “Bias Reduction of Maximum Likelihood Estimates”. In: *Biometrika* 80.1, pp. 27–38.

- Fishkin, James S. (2009). *When the People Speak: Deliberative Democracy & Public Consultation*. New York: Oxford University Press.
- Fong, Christian and Justin Grimmer (forthcoming). “Causal Inference with Latent Treatments”. In: *American Journal of Political Science*.
- Gaertner, Samuel L. and John P. McLaughlin (Mar. 1983). “Racial Stereotypes: Associations and Ascriptions of Positive and Negative Characteristics”. In: *Social Psychology Quarterly* 46.1, pp. 23–30.
- Gaines, Brian J. et al. (Nov. 2007). “Same Facts, Different Interpretations: Partisan Motivation and Opinion on Iraq”. In: *The Journal of Politics* 69.4, pp. 957–974.
- Gerber, Marlène et al. (Sept. 2014). “Deliberative and Non-Deliberative Persuasion: Mechanisms of Opinion Formation in EuroPolis”. In: *European Union Politics* 15.3, pp. 410–429.
- Gibson, James L. (July 1998). “A Sober Second Thought: An Experiment in Persuading Russians to Tolerate”. In: *American Journal of Political Science* 42.3, p. 819. JSTOR: 2991731.
- Gretz, Shai et al. (Apr. 2020). “A Large-Scale Dataset for Argument Quality Ranking: Construction and Analysis”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 34.5, pp. 7805–7813.
- Guess, Andrew and Alexander Coppock (Oct. 2020). “Does Counter-Attitudinal Information Cause Backlash? Results from Three Large Survey Experiments”. In: *British Journal of Political Science* 50.4, pp. 1497–1515.
- Habernal, Ivan and Iryna Gurevych (Apr. 2017). “Argumentation Mining in User-Generated Web Discourse”. In: *Computational Linguistics* 43.1, pp. 125–179.
- Hahn, Ulrike and Mike Oaksford (July 2007). “The Rationality of Informal Argumentation: A Bayesian Approach to Reasoning Fallacies”. In: *Psychological Review* 114.3, pp. 704–732.
- Haselton, Martie G., Daniel Nettle, and Paul W. Andrews (2005). “The Evolution of Cognitive Bias”. In: *The Handbook of Evolutionary Psychology*. Ed. by David M. Buss. 1st ed. Hoboken, NJ: John Wiley & Sons, Inc., pp. 724–746.

- Heinze, Georg and Michael Schemper (Aug. 2002). “A Solution to the Problem of Separation in Logistic Regression”. In: *Statistics in Medicine* 21.16, pp. 2409–2419.
- Hopkins, Daniel J. (Mar. 2014). “The Consequences of Broader Media Choice: Evidence from the Expansion of Fox News”. In: *Quarterly Journal of Political Science* 9.1, pp. 115–135.
- Itzchakov, Guy, Avraham N. Kluger, and Dotan R. Castro (Jan. 2017). “I Am Aware of My Inconsistencies but Can Tolerate Them: The Effect of High Quality Listening on Speakers’ Attitude Ambivalence”. In: *Personality and Social Psychology Bulletin* 43.1, pp. 105–120.
- Iyengar, Shanto and Kyu S. Hahn (Mar. 2009). “Red Media, Blue Media: Evidence of Ideological Selectivity in Media Use”. In: *Journal of Communication* 59.1, pp. 19–39.
- Jackman, Simon and Paul M. Sniderman (May 2006). “The Limits of Deliberative Discussion: A Model of Everyday Political Arguments”. In: *The Journal of Politics* 68.2, pp. 272–283.
- Jost, John T. et al. (May 2003). “Political Conservatism as Motivated Social Cognition”. In: *Psychological Bulletin* 129.3, pp. 339–375.
- Kahneman, Daniel (2011). *Thinking, Fast and Slow*. New York: Farrar, Straus and Giroux.
- Kahneman, Daniel, Paul Slovic, and Amos Tversky (1982). *Judgment Under Uncertainty: Heuristics and Biases*. Cambridge, UK: Cambridge University Press.
- Kalla, Joshua L. and David E. Broockman (June 2020). “Reducing Exclusionary Attitudes through Interpersonal Conversation: Evidence from Three Field Experiments”. In: *American Political Science Review* 114.2, pp. 410–425.
- Kinder, Donald R. and Nathan P. Kalmoe (2017). *Neither Liberal nor Conservative: Ideological Innocence in the American Public*. Chicago: The University of Chicago Press.
- King, Gary and Langche Zeng (2001). “Logistic Regression in Rare Events Data”. In: *Political Analysis* 9.2, pp. 137–163.
- Knox, Dean, Christopher Lucas, and Wendy K. Tam Cho (2022). “Testing Causal Theories with Learned Proxies”. In: *Annual Review of Political Science* 25, pp. 419–441.
- Koenig, Melissa A. (May 2012). “Beyond Semantic Accuracy: Preschoolers Evaluate a Speaker’s Reasons”. In: *Child Development* 83.3, pp. 1051–1063.



- Kuhn, Deanna (Sum. 1992). "Thinking as Argument". In: *Harvard Educational Review* 62.2, pp. 155–178.
- Kunda, Ziva (Nov. 1990). "The Case for Motivated Reasoning". In: *Psychological Bulletin* 108.3, pp. 480–498.
- Lebo, Matthew J. and Daniel Cassino (Dec. 2007). "The Aggregated Consequences of Motivated Reasoning and the Dynamics of Partisan Presidential Approval". In: *Political Psychology* 28.6, pp. 719–746.
- Levendusky, Matthew S., James N. Druckman, and Audrey McLain (Apr. 2016). "How Group Discussions Create Strong Attitudes and Strong Partisans". In: *Research & Politics* 3.2, pp. 1–6.
- Lodge, Milton and Charles Taber (2000). "Three Steps toward a Theory of Motivated Political Reasoning". In: *Elements of Reason: Cognition, Choice, and the Bounds of Rationality*. Ed. by Arthur Lupia, Mathew D. McCubbins, and Samuel L. Popkin. New York: Cambridge University Press, pp. 183–213.
- Lukin, Stephanie M. et al. (Aug. 2017). "Argument Strength Is in the Eye of the Beholder: Audience Effects in Persuasion". Pre-Print. University of California, Santa Cruz. arXiv: 1708.09085.
- Lupia, Arthur (Sum. 2002). "Deliberation Disconnected: What It Takes to Improve Civic Competence". In: *Law and Contemporary Problems* 65.3, p. 133. JSTOR: 1192406.
- Lupia, Arthur and Mathew D. McCubbins (1998). *The Democratic Dilemma: Can Citizens Learn What They Need to Know?* New York: Cambridge University Press.
- Luskin, Robert C., James S. Fishkin, and Roger Jowell (July 2002). "Considered Opinions: Deliberative Polling in Britain". In: *British Journal of Political Science* 32.3, pp. 455–487.
- Luskin, Robert C., Ian O'Flynn, et al. (Mar. 2014). "Deliberating across Deep Divides". In: *Political Studies* 62.1, pp. 116–135.
- Madison, James (Nov. 1787). "Federalist No. 10: The Same Subject Continued: The Union as a Safeguard Against Domestic Faction and Insurrection". In: *New York Daily Advertiser*.

- Madison, James (Feb. 1788). “Federalist No. 57: The Alleged Tendency of the New Plan to Elevate the Few at the Expense of the Many”. In: *New York Daily Advertiser*.
- Meckstroth, Christopher (Aug. 2012). “Socratic Method and Political Science”. In: *American Political Science Review* 106.3, pp. 644–660.
- Mellers, Barbara et al. (Mar. 2015). “The Psychology of Intelligence Analysis: Drivers of Prediction Accuracy in World Politics”. In: *Journal of Experimental Psychology: Applied* 21.1, pp. 1–14.
- Mendelberg, Tali and John Oleske (Apr. 2000). “Race and Public Deliberation”. In: *Political Communication* 17.2, pp. 169–191.
- Mercier, Hugo and Hélène Landemore (Apr. 2012). “Reasoning Is for Arguing: Understanding the Successes and Failures of Deliberation”. In: *Political Psychology* 33.2, pp. 243–258.
- Mercier, Hugo and Dan Sperber (Apr. 2011). “Why Do Humans Reason? Arguments for an Argumentative Theory”. In: *Behavioral and Brain Sciences* 34.2, pp. 57–74.
- (2017). *The Enigma of Reason*. Cambridge, MA: Harvard University Press.
- Munro, Geoffrey D., Carrie Weih, and Jeffrey Tsai (May 2010). “Motivated Suspicion: Asymmetrical Attributions of the Behavior of Political Ingroup and Outgroup Members”. In: *Basic and Applied Social Psychology* 32.2, pp. 173–184.
- Myers, David G. and George D. Bishop (Aug. 1970). “Discussion Effects on Racial Attitudes”. In: *Science* 169.3947, pp. 778–779.
- Parsons, Bryan M. (June 2010). “Social Networks and the Affective Impact of Political Disagreement”. In: *Political Behavior* 32.2, pp. 181–204.
- Payne, John W., James R. Bettman, and Eric J. Johnson (July 1988). “Adaptive Strategy Selection in Decision-Making”. In: *Journal of Experimental Psychology: Learning, Memory, and Cognition* 14.3, pp. 534–552.
- Pew Research Center (Apr. 7, 2021). *Social Media Fact Sheet*. URL: <https://www.pewresearch.org/internet/fact-sheet/social-media/?menuItem=3814afe3-3f3c-4623-910b-8a6a37885ab8>.

- Prior, Markus (May 2013). "Media and Political Polarization". In: *Annual Review of Political Science* 16.1, pp. 101–127.
- Puhr, Rainer et al. (2017). "Firth's Logistic Regression with Rare Events: Accurate Effect Estimates and Predictions?" In: *Statistics in Medicine* 36, pp. 2302–2317.
- Redlawsk, David P. (Nov. 2002). "Hot Cognition or Cool Consideration? Testing the Effects of Motivated Reasoning on Political Decision Making". In: *The Journal of Politics* 64.4, pp. 1021–1044.
- Sandoval, William A. (Jan. 2003). "Conceptual and Epistemic Aspects of Students' Scientific Explanations". In: *Journal of the Learning Sciences* 12.1, pp. 5–51.
- Schneiderhan, Erik and Shamus Khan (Mar. 2008). "Reasons and Inclusion: The Foundation of Deliberation". In: *Sociological Theory* 26.1, pp. 1–24.
- Settle, Jaime E. (2018). *Frenemies: How Social Media Polarizes America*. New York: Cambridge University Press.
- Sniderman, Paul M., Richard A. Brody, and Philip E. Tetlock (1991). *Reasoning and Choice: Explorations in Political Psychology*. New York: Cambridge University Press.
- Sniderman, Paul M. and Thomas Piazza (1993). *The Scar of Race*. Cambridge, MA: Belknap Press.
- Stanley, Matthew L., Paul Henne, et al. (Sept. 2020). "Resistance to Position Change, Motivated Reasoning, and Polarization". In: *Political Behavior* 42.3, pp. 891–913.
- Stanley, Matthew L., Peter S. Whitehead, et al. (Nov. 2020). "Exposure to Opposing Reasons Reduces Negative Impressions of Ideological Opponents". In: *Journal of Experimental Social Psychology* 91.
- Stigler, George J. (June 1961). "The Economics of Information". In: *Journal of Political Economy* 69.3, pp. 213–225.
- Sunstein, Cass R. (June 2002). "The Law of Group Polarization". In: *The Journal of Political Philosophy* 10.2, pp. 175–195.
- Taber, Charles S., Damon Cann, and Simona Kucsova (June 2009). "The Motivated Processing of Political Arguments". In: *Political Behavior* 31.2, pp. 137–155.

- Taber, Charles S. and Milton Lodge (July 2006). “Motivated Skepticism in the Evaluation of Political Beliefs”. In: *American Journal of Political Science* 50.3, pp. 755–769.
- (Feb. 2016). “The Illusion of Choice in Democratic Politics: The Unconscious Impact of Motivated Political Reasoning”. In: *Political Psychology* 37 (Suppl. 1), pp. 61–85.
- Tesler, Michael (Oct. 2015). “Priming Predispositions and Changing Policy Positions: An Account of When Mass Opinion Is Primed or Changed”. In: *American Journal of Political Science* 59.4, pp. 806–824.
- Tetlock, Philip E. and Dan Gardner (2015). *Superforecasting: The Art and Science of Prediction*. New York: Crown Publishers.
- Trouche, Emmanuel, Petter Johansson, et al. (Nov. 2016). “The Selective Laziness of Reasoning”. In: *Cognitive Science* 40.8, pp. 2122–2136.
- Trouche, Emmanuel, Emmanuel Sander, and Hugo Mercier (Oct. 2014). “Arguments, More than Confidence, Explain the Good Performance of Reasoning Groups”. In: *Journal of Experimental Psychology: General* 143.5, pp. 1958–1971.
- VanderWeele, Tyler J. and Miguel A. Hernán (June 2012). “Results on Differential and Dependent Measurement Error of the Exposure and the Outcome Using Signed Directed Acyclic Graphs”. In: *American Journal of Epidemiology* 175.12, pp. 1303–1310.
- Walker, Marilyn A. et al. (May 2012). “A Corpus for Research on Deliberation and Debate”. In: *Proceedings of the Eighth International Conference on Language Resources and Evaluation*. Istanbul, pp. 812–817.
- Wang, Lu and Claire Cardie (June 2014). “Improving Agreement and Disagreement Identification in Online Discussions with A Socially-Tuned Sentiment Lexicon”. In: *Proceedings of the 5th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*. Baltimore: Association for Computational Linguistics, pp. 97–106.
- Westwood, Sean J. (Oct. 2015). “The Role of Persuasion in Deliberative Opinion Change”. In: *Political Communication* 32.4, pp. 509–528.

- Wood, Thomas and Ethan Porter (Mar. 2019). “The Elusive Backfire Effect: Mass Attitudes’ Steadfast Factual Adherence”. In: *Political Behavior* 41.1, pp. 135–163.
- Wooldridge, Jeffrey M. (2015). *Introductory Econometrics: A Modern Approach*. Mason, OH: Cengage Learning.
- Zaller, John and Stanley Feldman (Aug. 1992). “A Simple Theory of the Survey Response: Answering Questions versus Revealing Preferences”. In: *American Journal of Political Science* 36.3, p. 579.
- Zaller, John R. (1992). *The Nature and Origins of Mass Opinion*. New York: Cambridge University Press.
- Zhu, Yukun et al. (Dec. 2015). “Aligning Books and Movies: Towards Story-Like Visual Explanations by Watching Movies and Reading Books”. In: *2015 IEEE International Conference on Computer Vision*. Santiago, Chile: Institute of Electrical and Electronics Engineers, pp. 19–27.