

Critical Thinking and Storytelling

Brian Jabarian¹ and Elia Sartori²

¹Booth Business School, University of Chicago

²Center for Studies in Economics and Finance and Università degli Studi di Napoli Federico II.

First version: September 5, 2020;
Current version: October 18, 2023*

Abstract

In a lab-in-field online experiment on a representative US population ($N = 725$), incentivized through an LLM, we show that different storytelling formats – different media sources and styles presenting the same set of facts – affect the intensity at which individuals become critical thinkers, that is, become aware of a trade-off between competing worldviews. Intermediate storytelling formats (Facebook) are more effective in triggering individuals to think critically than shorter basic storytelling formats (Twitter) and more extended and sophisticated storytelling formats (newspaper). Individuals with a high need for cognition drive the differential effects of treatments, underscoring the importance of cognitive styles in storytelling personalization for the digital content economy. Finally, in a stylized voting model, we explore the role of critical thinking, cognitive personalization, and storytelling formats in digital voting contexts. We establish that increasing the share of *critical thinkers* in the population increases the efficiency of surveys and elections but might increase the bias of elections (or surveys).

*Contact: brian.jabarian@chicagobooth.edu

We are indebted to Roland Bénabou for his continued guidance. We are grateful for the comments of Tore Ellingsen, Nicolas Jacquemet, Yves Le Yaouanq, John List, Dan McGee, Pietro Ortoleva, Devin Pope, Eldar Shafir, Avner Strulov-Shlain, Richard Thaler, Jean-Marc Tallon, Marie-Claire Villeval, George Wu, Leeat Yariv, and Sam Zbarsky. The authors are grateful for the research assistance provided by Christian Kontz, Andras Molnar, Alessandro Sciacchitano, and Alfio De Angelis. We are grateful for the participation of psychologists from the Department of Psychology at Princeton University. We thank the seminar attendees at Bologna, PSE. We obtained Princeton IRB approval #12995 on June 12, 2020. This paper was written in part while Brian visited the Department of Economics at Princeton University and the Kahneman-Treisman Center for Behavioral Science and Public Policy in 2018-2020. He thanks their hospitality. He also acknowledges financial support from the Paris School of Economics, the Sorbonne Economics Center, and the Forethought Foundation for his visit to Princeton, Grant ANR-17-CE26-0003, and Grant ANR-17-EURE-001. *Full Acknowledgment to be added.*

1 Introduction

Individuals rely on a variety of reasoning styles to form preferences on trade-offs or dilemmas. These latter co-defendants of views and policies, without reaching a definitive conclusion, trigger ambivalent attitudes within agents [Kaplan \(1972\)](#), leading to the formation of their preferences. This formation can happen through different mental models: in a fast, intuitive way or slow, reasoned way (Kahneman). Within reasoned mental models, agents can navigate between naive thinking and critical thinking. A key difference between both styles lies in the awareness that the issue at hand is a trade-off or dilemma ([Halpern \(2013\)](#)).

This paper presents a novel and simple incentivized experiment to identify and classify transitions between naive thinking and critical thinking. Our design shows how to use storytelling formats to identify such a transition in mental models and not only, as often shown in the literature, a shift in preferences. Consequently, this research contributes to the growing experimental literature that aims to identify such critical thinking and examine their impact on policy ([List \(2022\)](#)). It further expands the spectrum of reasoning styles investigated in the behavioral literature, which is usually focused on motivated reasoning ([Kunda \(1990\)](#), [Bénabou and Tirole \(2006\)](#)).

In addition, our study emphasizes the role of tailored storytelling formats, where factual information is conveyed through a specific visual design and writing style – commonly known as “UX design” in the marketing and communication literature. This is instrumental in shifting individuals from stereotypical to critical thinking when hard facts are absent. The interplay of quantity, quality, and personal cognitive styles of information is crucial in shaping preferences. In the digital economy, the media acts as a significant catalyst, encouraging critical thinking. When confronted with dilemmas, individuals find that mere reliance on ‘objective facts’ is insufficient. Instead, critical thinking is key to developing well-reasoned preferences.

Consequently, we adopt the terms “story” and “storytelling formats” to represent “media content” and “media format,” respectively. The manner in which an issue is presented, ranging from a simplistic tweetstorm to a detailed newspaper article, can influence individual awareness of an issue’s ambivalence. Our main findings indicate that two similar interpretations of the same fact, presented through different visual formats and writing styles, elicit different behavioral responses in participants. This expands the conventional definition of “narratives” in economics, generally defined as a specific interpretation of a fact ([Shiller \(2017\)](#), [Eliaz and Spiegler \(2020\)](#)), by considering the specific format through which the interpretation is presented. In summary,

when it comes to dilemmas, the influence of information on individual preferences lies not just in the content, but also in the delivery.

Our experimental design involves three primary stages. Initially, we categorize participants as stereotypical or critical thinkers on a contentious topic using a combination of self-report measures and incentivized elicitation techniques. Subsequently, we expose participants to one of three storytelling interventions, each centered on the same set of pros and cons related to the issue. The storytelling strategies range from a concise, simplistic style presented in a Twitter-like format to a medium-level complexity style in a Facebook-like format, to an intricate, detailed style in a newspaper-like format.

In considering social media through the lens of these storytelling perspectives, we contribute to a growing body of literature that increasingly focuses on the impact of specific *formats* on shaping individual political behaviors, such as voting (Gorodnichenko et al. (2021), Munir (2018), Falck et al. (2014)). More generally, our findings highlight an additional channel (altering the share of critical thinkers) through which social networks can affect welfare, contributing to the rapidly expanding literature on social networks and welfare (Allcott et al. (2020)).

After exposure, participants are prompted to write an incentivized critical thinking essay following specific guidelines. Completing this task is incentivized using Large Language Models (LLMs), particularly GPT-3, which offers an automatic comparative ranking against a US average score. Given the absence of a definitive “critical thinking” measure, we instituted a secondary experiment to collect expert human feedback. This practice aligns with the model-based reinforcement learning techniques commonly employed by AI-oriented companies, such as OpenAI. Expert feedback was obtained from cognitive psychologists with Ph.D. or higher degrees and experience in ambivalence and critical thinking. The essays authored by the participants were randomly assigned to three independent expert labelers tasked with grading the submissions as pass or fail, depending on their judgment of clear indications of critical thinking in the content. The data from these evaluations were then used to reclassify the participants as stereotypical or critical thinkers after the storytelling interventions.

Using psychological literature on cognitive sophistication, we measured participants’ need for cognition (Cacioppo and Petty (1982)) and cognitive flexibility (Martin and Rubin (1995)) throughout our experiment. These metrics enabled us to conduct a heterogeneity analysis. Intriguingly, our results indicate that people with a higher need for cognition transition more quickly from stereotypical to authentic preferences upon exposure to a medium level of storytelling (i.e., Facebook) than a lower (Twitter)

or higher one (Newspaper).

Upon establishing the role of storytelling as a catalyst for critical thinking, we examined its implications for industrial organization and political economy. Specifically, we explore how storytelling techniques impact the efficiency of surveys and elections, making them crucial to social welfare in industrial and political decision-making contexts. In fact, decision-makers who use surveys and run elections may find that critical thinking preferences offer more reliable data than raw preferences.

Consider a public figure or organization whose social image or economic returns hinge on the public endorsement of their position on a particular issue. Such a principal needs to anticipate the public's expected stance, as public endorsements serve as reputational commitments and "focusing events" that prompt individuals to evaluate their raw preferences and establish reasoned preferences critically. Therefore, the principal should gauge the public's reasoned preferences before declaring a stance, minimizing the risk of sustained backlash. Suppose that such an estimate is based on a poll. In that case, its precision depends on the respondents reporting their reasoned preferences, which requires that these preferences have been formed in the first place.¹

Furthermore, consider an institutional principal, such as a policymaker, tasked with formulating an economic policy on a societal issue that presents a binary dilemma. The principal can select from a wide spectrum of policy alternatives. The optimal policy is a function of the distribution of reasoned preferences, i.e., the proportion of individuals who prefer one alternative over the other after engaging in critical thinking. This establishes the need for the principal to anticipate (and incentivize) the formation of agents' reasoned preferences before making a decision, since the reasoned preference distribution forms its normative criterion.²

In both instances, we identify a principal who is interested in the distribution of reasoned preferences: either out of fear that their actions will provoke a backlash if they deviate excessively from the target or because they use such a distribution "for lack of anything better" as an appropriate normative criterion for the social aggregation of preferences. Elections would be efficient if all individuals reported their reasoned preferences at the poll, allowing a precise estimation of the relevant unknown. However, individuals arrive at their reasoned preferences only after participating in a critical

¹It is conceivable that in a strategic voting setting, agents might misreport their reasoned preferences even after forming one. However, we view this concern as secondary to our intended application. Hence, we assume that the formation and reporting of a reasoned preference are congruous actions.

²Ultimately, such policymakers must adopt a policy aligned with one of two conflicting worldviews. Most of the time, if they were critical thinkers, they would recognize their rational preference on the issue but risk imposing it on the rest of the population.

thinking process, a process that not all individuals may have completed by the time the election is held.³

The remainder of this paper is organized as follows. Section 2 details our experimental design. Section 3 discusses our empirical findings. Section 4 presents our behavioral model, accompanied by its key positive and normative results. Finally, Section 5 provides a conclusion in which we discuss potential limitations and future extensions of our model and experiment.

Predictive Power of Elections. As demonstrated in the seminal work of Feddersen and Pesendorfer (1997), there are instances where many voters effectively aggregate information, resulting in an equilibrium outcome that is fully information equivalent. However, preference heterogeneity can impede a voting procedure from effectively aggregating individual voters' information (Kim and Fey (2007); Gul and Pesendorfer (2009); Bhattacharya, 2013; Acharya, 2016; Ali et al., 2018). This literature highlights that when voters have divergent preferences and incomplete information about the state of nature, they may collectively choose an outcome that is less favorable for society or preferred by a minority. The central concern in these papers is strategic voting, an issue not present in our analysis. In our model, every citizen votes for their current preference. Still, the extent to which this preference accurately reflects the payoff-relevant reasoned preference depends on the citizen's cognitive state, which is influenced by politics. Political can be seen as a method of making citizens view their reasoned preference as private information that an election aims to uncover.

Social Media and Welfare. Given the rapid growth of social networks, the literature has started to focus on its impact on several economic variables. Allcott et al. (2020) show by means of a randomized experiment that social networks (Facebook) undermine the welfare of the individual. The researchers used a randomized experiment to evaluate the impact of Facebook on the welfare of individuals and found that social media undermined the welfare of agents.

Although social networks seem to harm the welfare, the results have little meaning. The point that we strive to make is that the way news is presented matters in making people realize ambivalence about an issue. It could well be that all social networks reduce the attention to experience and therefore "stuckpeople" (that is, it reduces λ and possibly skews the distribution of preferences away from the real mean μ_Y). Our

³We posit that the principal cannot "screen" voters based on their stage of critical thinking and thus only utilize "informed voters". However, our findings suggest that a survey methodology capable of categorizing agent types could significantly improve its precision.

analysis is necessarily partial and, to some degree, our experiment forced normative welfare into analytical thinking because it forced them to write an essay.

Not only does information and social media impact voting behavior, but the literature shows that social media can also influence another kind of behavior. [Gentzkow and Shapiro \(2010\)](#) showed that the Internet changes the ideological segregation among American voters and highlights heterogeneity depending on online and offline news. [?](#) explored the causal effect of social networks on hate crime, specifically anti-refugee sentiment. They found that hate crime was more notable in municipalities with higher levels of social media use.

Structure of the paper. Section 4 describes the behavioral model and its main positive and normative results. Section 2 details the experimental design. Section 3 elaborates on the empirical results. Section 5 concludes and discusses possible caveats and extensions of the model and experiment.

2 Design

2.1 Overview

In a nutshell, in our experiment, we expose subject participants to different storytelling formats and elicit their pre- and post-treatment stages in the critical thinking process associated with a dilemma. We then test whether the likelihood of transitioning from stereotypical thinkers, S , to critical thinkers, A , varies significantly between formats. Throughout the experiment, we also collected data about participants' cognitive styles using standard measures from the psychological literature. This allows us to test whether the effectiveness of certain storytelling formats is achieved through identifiable cognitive traits. We used incentivized elicitation for key individual variables -pre and post-awareness states - and implemented anti-cheating policies and attention screeners to ensure optimal data collection quality.⁴ Figure 1 provides an overview of the experimental design and its primary elicitation, which we will elaborate on in subsequent sections.

We gathered 900 participants from a representative US population using Prolific, a data collection platform increasingly favored by economists due to its high data quality. Following a meticulous screening for attention, cheating and quality, as outlined

⁴The Princeton Institutional Review Board approved the experiment. See the appendix for the detailed Princeton IRB approval.

in the following sections, our final sample size was $N = 725$. Participants received a fixed payment of \$2 and a bonus payment of up to \$5, resulting in an average payment of approximately \$6.

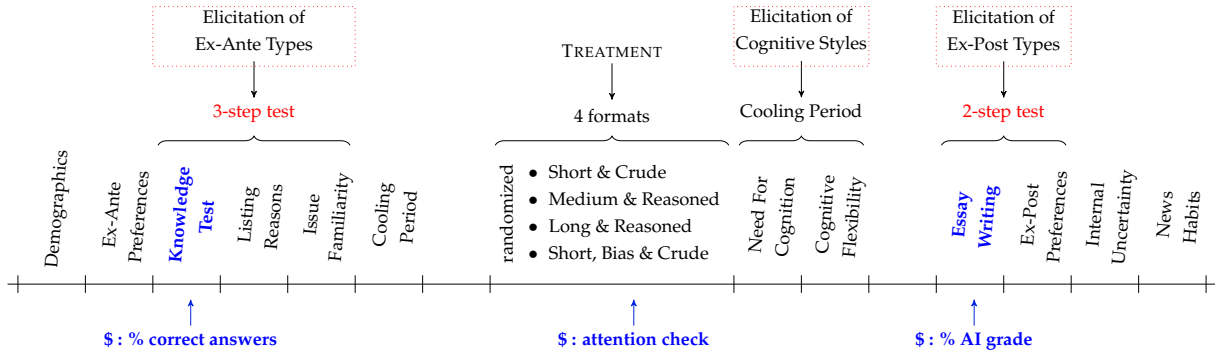


Figure 1: Experiment Design

2.2 Classification of Mental Models

We now present and describe the rationale for the two strategies we employ to classify participants into reasoning states $\{S, A\}$ before and after treatments. Table 1 summarizes both strategies.

Pre-treatment classification strategy. We use a three-pronged test to design our pre-treatment classification strategy as $\{S, A\}$. This test is based on the following heuristic conditions critical thinkers must satisfy: i) they must have basic knowledge of the issue at hand; ii) they must have thought about the issue before; iii) they must be aware that there exist both pros and cons for the issue. All i)-iii) characteristics are needed to be a critical thinker about an issue to avoid misclassification (as it could be by only using iii)).

To generate condition i), we rely on an assessment designed by Pew Research (Vogels and Anderson (2019)) and launched on a representative US population, referred to as the *knowledge test* in Figure 1. In our experiment, to pass the knowledge test, participants must score at least as high or higher than the nationally representative US

population score found by Pew Research.⁵ To elicit condition ii), we ask participants to report whether they have thought about the issue before coming to our experiment. To elicit Condition iii), we ask them to provide evidence by providing two reasons that support their preference for the digital privacy issue and two that go against their preference. This task is instrumental in more accurately targeting the state of awareness of the state of prior treatment of individuals. We refer to this task as *The list of reasons* is Figure 1. If (and only if) subjects are already beyond their raw preference stage, we can provide a complete classification.

We cannot rely on the same three-pronged tests to provide the post-treatment classification of the participants in terms of $\{S, A\}$. Since the pretreatment classification test includes condition iii) and our storytelling format treatments expose subjects to a series of pros and cons about the issue (see the next section), relying on the same condition here can misidentify critical thinking as memory effects. Indeed, subjects might not be critical thinkers, having accepted the issue as ambivalent by default, but happen to remember their list of pros and cons reported before treatment. Therefore, we need a different post-treatment classification strategy.

Post-treatment classification strategy. We classify participants' post-treatment critical thinking state as follows. We require participants to write an incentivized essay discussing their preferences on the issue at hand. Subjects are instructed to present the issue and articulate their argumentative position⁶ Their payment is based on the quality of their article measured by a software powered by a large language model (generative AI), Grammarly.

Although Grammarly is efficient in assessing the overall quality of writing (at the time of our experiment, still powered by models similar to GPT-3), it lacks the capacity to capture the nuances of critical thinking, especially in terms of discerning whether the writer demonstrates an awareness of the ambivalence surrounding the issue. To address this limitation, we ask cognitive psychologists with Ph.D. degrees to provide a professional assessment of the essays. These experts are randomly assigned to the participants' essays and are asked to evaluate whether the essay reflects a state of awareness or not, assigning a pass or fail grade accordingly. While participants receive payment based on the AI's evaluation, our analysis focuses on the cognitive psychologists' assessment, with AI's scores serving as a robustness check (see Section 3.3).

⁵It consists of 10 questions. See the appendix for the wording details.

⁶See Appendix to see the specific instructions to participants.

Treatment	A	S
BEFORE	Knowledge Test Score $> \tau_{KTS}$	
	Issue Familiarity = 1	Else
	Reasons List $> \tau_{RL}$	
AFTER	Psychologists Grade = Pass	Else

Table 1: CLASSIFICATION STRATEGY BEFORE / AFTER TREATMENT

2.3 Measuring Cognitive Styles

Since we are interested in explaining possible drivers of our results, through the experiment, we measure participants’ cognitive styles and correlate those with the treatments effectiveness. We measure cognitive styles according to three metrics. The first two are standard in the psychology literature: the Need for Cognitive Scale (NCS) and the Cognitive Flexibility Scale (CFS).

NCS measures a participant’s willingness to think deeply. It was proposed by neuroscientists and cognitive psychologists [Cacioppo and Petty \(1982\)](#) and has become a gold standard in cognitive psychology. It comprises a series of six questions that each receive a score between 1 and 5. We compare the aggregate score with the sample average to classify participants into high or low need for cognition.⁷

CFS measures an agent’s ability to switch between thoughts and courses of action. It was proposed by cognitive psychologist [Martin and Rubin \(1995\)](#) and is a standard scale in cognitive psychology. It comprises a series of six questions that each receive a score between 1 and 6. We compare the aggregate score with the average of the US population to classify participants as high or low cognitive flexibility.⁸

Finally, we use the AI generated score of an essay, unrelated to the core issue of our experiment, measuring an individual’s abilities to coherently present an argument.

⁷No average for the US population is available for this score, despite being used widely across the social sciences and psychology. In addition, it was originally developed as a 34-question version, but the authors developed a shorter, more efficient 18-question version to elicit other psychological characteristics during the same laboratory session. Since then, it has been considered the benchmark scale widely used in the cognitive and social sciences. An even shorter 6-question version has been tested and validated, allowing it to be implemented in a field survey experiment in which the participant’s attention is even more scarce. We will use this later.

⁸The average is provided by the authors: 55.

2.4 Description of Treatments

Participants are randomly assigned to one of four treatments. These treatments contain the same content (that is, the same selection of facts about the digital issue, but differ in semantic style and graphic design, as elaborated in the Introduction) and last the same amount of time.

In summary, such formats range from the semantically crudest presentation of facts to the most refined presentation. The TWITTER treatment presents them more crudely through a “tweet” format. The FACEBOOK treatment uses the format of “Facebook posts.” The NEWSPAPER treatment presents them in the most refined way through “newspaper articles.” The PARTISAN TWITTER treatment uses only a partisan Twitter format (either only pros or cons)⁹.

Before treatment starts, participants are explicitly informed that despite their high similarity to real news, Facebook tweets and posts are fake. At the end of the experiment, participants were briefed and reminded that tweets and Facebook posts were fake, following common practice in behavioral and experimental economics and according to our Institutional Review Board (IRB) approval.

2.5 Incentive Mechanisms and Quality Screening

Incentive mechanism. In the experiment, participants receive two types of payment. First, they receive a fixed reward of \$2 for fully completing the experiment by answering the comprehension questions correctly, guaranteed. Second, they receive a bonus payment, at most \$6, as described below.

Most of the participants’ bonus payments (up to 5\$; participants’ performance in the writing exercise, which captures their critical thinking process, determines their bonus). We ask participants to write two short essays during this study that will be graded from 0 to 100 points using Artificial Intelligence (AI)-powered Grammarly software¹⁰ We divide the bonus payment into two parts.

The largest part (from \$0 to \$5) is proportional to the weighted average score on the essay writing task; the second essay receives more weight (2/3) because it requires more writing (400 characters as opposed to 200 characters). The score can range from 0 to 100 points and the reward will be proportional to the score. If the participants score 0, they win \$0. If they get a score of 50, they win \$2.50. If they get a score of 100,

⁹In the appendix we detail and provide examples of each treatment.

¹⁰We, the authors, confirm to have neither professional ties nor a business contract with this company. See the appendix for a summary of how this AI works.

then they win \$5. An essay that receives a low score from the AI can still earn a high score on critical thinking and awareness, despite the writer’s difficulty with English. In the instructions to psychologist graders, we define and exemplify what we mean by a “dilemma,” “realizing that the issue is ambivalent,” and “critical thinking”. We also run robust checks with philosophers.

To be eligible for the remaining bonus payment (up to \$1), participants must receive at least an average score of 50/100 in the essay exercise in addition to the bonus of the writing essay. This requirement ensures that participants take the exercise seriously; cheaters and agents that are inconsistent in their preferences are not eligible for this bonus payment. The performance of the participant on the knowledge test determines this additional bonus. The test consists of 10 questions and each participant receives \$0.10 for each correctly answered question.

Monitoring Algorithms for Cheating Behavior. We implement three attention screeners as is standard in online experimental economics. The core of our experiment is for participants to write an original essay by themselves. We need the subjects to avoid accessing external information during the writing task. As such, we implement two algorithms to monitor cheating behavior.

Before starting their experiment and on par with the IRB, we inform participants that they must not access external information during the experiment, particularly during the knowledge test and essay exercise. In addition, the essay must be original. Failing to do so would be considered “cheating behavior.” As such, they would be red-flagged and prevented from receiving anything other than the fixed payment. We excluded such participants from our data analysis.¹¹

The first algorithm tracks the number of times that the participants open a new tab on their computer during the essay exercise and how much time they spend on our essay writing web-page¹². The second algorithm checks whether participants copy-paste external information by comparing the number of written characters and the number of keyboard clicks. If the number of keyboard clicks is strictly inferior to the number of written characters, this implies that the participants have copied external information. This second algorithm cannot distinguish between the original external information¹³ and plagiarism. Therefore, we use a feature in the AI software to check

¹¹We provide both algorithms as open source in our [GitHub](#).

¹²For legal privacy purposes, we did not access the content of the opened tab but gathered only the following information: ‘participant i has opened a new tab during the essay, n number of times, for such and such period t .

¹³In the situation in which some participants had already written on the topic or a relevant topic and

for plagiarism after the participants have finished the experiment.

3 Analysis

3.1 Storytelling Prompts Critical Thinking

We now test whether storytelling formats have a role in the critical thinking process of individuals. To this end, we calculate, for each treatment $i = \{newspaper, twitter, facebook\}$, the frequency $\hat{\lambda}_i$ with which agents subject to the format i transition from the critical thinking state S to the critical thinking state A . Formally,

$$\hat{\lambda}_i = \frac{\#(S \rightarrow A)_i}{\#(S \rightarrow A)_i + \#(S \rightarrow S)_i}$$

We used estimated intensities to perform a difference-in-means test of the null hypotheses $\lambda_i = \lambda_j$ for all possible combinations of treatments $\{i, j\}$.

Table 2 collects point estimates and confidence intervals. From this table we observe that the only significant difference is between *Facebook* and *Twitter*, where the former performs better in transitioning subjects from critical thinking state S to A . Through this significant result, we establish that the format affects the critical thinking process. When exposed to a different way of presenting the same basic information, people realize the ambivalent nature of the issue at hand differently. In section –, we perform robustness checks of this result (different thresholds, etc., metrics for success) to understand potential drivers of this effect.

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.332 (0.054)	-0.865 (0.054)
TWITTER	.	.	-2.249** (0.053)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Table 2: z-score DIFFERENCE-IN-PROPORTIONS

saved it on their computer before coming to the experiment.

A possible explanation for the observed difference in the impact of storytelling formats on critical thinking, as highlighted in Table 2, is that treatment TWITTER may rely on a format that is too simplistic or naive to effectively push users toward critical thinking. This explanation can be elaborated on as follows.

First, while Twitter imposes a character limit on its content, forcing users to use concise language and simplifying complex ideas, Facebook allows for longer and more detailed posts¹⁴. This difference in content structure could affect how people process information and engage in critical thinking.

Second, the fast-paced nature of Twitter feeds and the emphasis on real-time information sharing could discourage users from pausing, reflecting, and analyzing the content they consume. This constant influx of new information might contribute to a deeper engagement with the material, reducing the likelihood of critical thinking.

Third, Twitter's focus on short, attention-grabbing headlines and sound bites may encourage users to form quick opinions based on surface-level information rather than delving deeper into the nuances of an issue. This aspect of the platform's design might hinder the development of well-informed perspectives and critical thinking.

Fourth, the prevalence of echo chambers on Twitter, where users primarily follow and interact with those who share their views, could further contribute to the observed limitations of the Twitter format in promoting critical thinking. This selective exposure to information might reinforce pre-existing beliefs and discourage users from challenging their assumptions.

Fifth, another factor to consider is the nature of user engagement on these platforms. Facebook is known for fostering more personal connections and allowing in-depth conversations, while Twitter primarily emphasizes short and quick information exchanges. This contrast in user engagement could contribute to the observed difference in the effectiveness of storytelling formats in critical thinking.

Finally, the role of media consumption habits might be influential in explaining the difference in critical thinking outcomes. Facebook users may be more inclined to read longer posts and engage in reflective thinking, whereas Twitter users may be more accustomed to quickly skimming through bite-sized information. As a result, individuals' media consumption habits could shape their receptiveness to the storytelling formats on these platforms, ultimately affecting their critical thinking process.

¹⁴This experiment was designed and launched before Musk Twitter's area, which led to the increase of tweets lengths for Blue Twitter users, which now can be considered as our Facebook treatment.

3.2 Cognitive Styles and Storytelling Personalization

We explore whether the cognitive traits we elicited explain the differential effect by conducting a split-sample difference in means. We test whether $\lambda_i = \lambda_j$ by partitioning our sample into high or low individuals in our cognitive metrics. The idea is that a more in-depth approach (such as the journal article) may be more effective for individuals more prone to think deeply. The efficacy of the *Facebook* treatment was driven by its differential impact on *High Need for Cognition*.

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	0.764 (0.070)	-2.238* (0.079)
TWITTER	.	.	-3.087** (0.075)
FACEBOOK	.	.	.

Standard errors in parentheses
 * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 3: z-score FOR *High Need for Cognition*.

These results suggest that the subjects who are most affected by the storytelling format are those who exhibit a high need for cognition. For them, treatment **FACEBOOK** seems to provide the right format to maximally capture their attention to present an issue so that it successfully nudges them to perform the critical thinking process.

One possible explanation for the results observed in Table 3 could be rooted in the characteristics of people with a high need for cognition. These individuals typically exhibit a greater tendency to engage in effort-based cognitive activities and prefer more complex information processing (Cacioppo and Petty (1982)). Consequently, the Facebook format could provide a more stimulating environment for critical thinking by offering a richer and more nuanced presentation of information than the Twitter format.

Furthermore, it has been suggested that individuals with a greater need for cognition are more likely to seek, attend to, and remember information consistent with their attitudes and beliefs (Hass & Linder, 1981). As a result, the Facebook format could be more effective in capturing your attention and motivating you to critically evaluate the content. This might explain why Facebook treatment significantly impacts transitioning subjects from critical thinking state S to A among those with a High Need for Cognition.

Future research could explore the specific features of the Facebook format that contribute to its efficacy in promoting critical thinking among individuals with a high need for cognition. For example, it would be interesting to investigate the role of multimedia elements, interactivity, and the integration of various information sources in fostering an environment conducive to critical thinking.

3.3 Robustness Checks

We address two potential challenges to ensure the robustness of our findings, threshold sensitivity and writing similarity checks, that we present now.

Threshold sensitivity. Our conclusions should remain consistent regardless of the specific values of the threshold used in characteristics i) and iii) of the three-pronged test we use to classify participants prior to treatment. Recall that i) refers to the digital knowledge test, and iii) refers to the reasons list exercise.

Regarding i), at the beginning of the study, we require the participants to score at least seven correct answers out of 10 questions. Compared to the original setting provided by Pew Research, our threshold is much more demanding. The Pew Research quiz was launched in a large representative sample in the US of 4,272 adults living in the United States. The median number of correct answers was four. Only 20% of the adults correctly answered seven or more questions and only 2% correctly answered the 10 questions. Despite this difference, we are still interested in determining whether our treatment effectiveness depends on scoring higher or lower than scoring 7 out of 10. Regarding iii), at the beginning of the study, we require participants to be able to list at least one reason for one side (pro or con) and two reasons for the other side (pro or con). We are interested in checking whether the effectiveness of our treatment depends on the ability of the participants to list more than one reason for each side.

Writing similarity. Our findings should not be influenced by the similarity in length between the essay task and any specific treatment, particularly the Facebook treatment. By comparing outcomes across different essay lengths or imposing length constraints, we can verify that the observed effects are not artifacts of such similarities, ensuring the robustness of our results.

In general, our robustness analysis confirms that the effectiveness of the treatment depends neither on the threshold sensitivity test nor on the writing similarity test.¹⁵

¹⁵We provide the analysis in the appendix.

3.4 Discussion of Empirical findings

Regarding the internal validity of our experiment, we recognize that classifying individuals' critical thinking states is inherently challenging. We devised different classification rules for pre- and post-treatment to avoid mistaking memory for critical thinking. Second, we use a noisy measure to look at the differences *between treatments*. Regarding the external validity of our experiment, the reader should refrain from interpreting our experiment as a comparison of social networks, concluding that "Facebook is better" but that rather "the format matters." In this interpretation, the entire class of social media becomes a storytelling format: one is exposed to a greater number of views, but they are possibly superficial. Does it help to become aware of the ambivalence of the issue about one's life experience or the in-depth study of a topic (more personal and reasoned, but time-consuming and unlikely to occur)? ¹⁶

Moreover, the analysis is also problem-specific. We have used digital privacy, but we realize that the nature of the problem might affect the features of the medium that make it more or less effective at inducing critical thinking. We realize that issues have an "objective" side where standard analysis of information acquisition from multiple sources is. However, we believe that a dimension of is inherent in many issues and that, beyond providing and helping to absorb "hard" information media, also have a role to propel individuals into realizing the ambivalent nature of the issue. It is along this dimension that our analysis has special relevance as it highlights a novel channel, possibly orthogonal to the ones identified before (Bernheim + [lit on social media]), through which. The same characteristics that make hard information difficult to absorb (e.g. continuous display with superficial language) might drive the attention on a specific issue and make individuals aware of its ambivalent nature. The theoretical model that we develop in the next section shows that those considerations are relevant for the efficiency of elections.

4 Theoretical Framework

In this section, we develop a simple model in which the intensity of the critical thinking process affects the efficiency of preference aggregation through elections. We consider a stylized social choice setting in which the utility is the distance between the political action and a target determined by the distribution of reasoned preferences, i.e. those held after completing the *critical thinking* process, within the population. This parame-

¹⁶extension of expert GPT-4 using our data set.

ter is unknown at first and can only be estimated using the result of a poll held at some time t , when (a part of) the citizens may still not have completed their process. Our main result — Propositions 1-*i*) — establish that such intensity is a relevant welfare measure: Regardless of the time of the elections, a higher intensity increases the information content of elections. Combining this result with the results of our experiment (Section 3.1) indicates that the way information is presented to agents (storytelling format) prior to “voting” affects the quality of information that is elicited in a poll.

We also show that the unambiguous comparative statics only holds if the principal can freely manipulate the results of the poll when taking her action (which we refer to as the P positive principal). If the election outcome constrained her action, as is most likely the case for an Institutional principal, then a bias-precision trade-off makes the comparative statics ambiguous: a faster critical thinking process might hurt efficiency of elections (Propositions 1-*ii*) and 2).

The section proceeds as follows. We first present the two welfare benchmarks (corresponding to the two types of principals discussed in the introduction) and the critical thinking process separately. Combining the two we then obtain closed-form (evolution of) welfare and establish under what conditions a faster critical thinking process is beneficial. Finally, we discuss the weakenings of our assumptions that seem most critical to building a richer model whose objective is not limited to establish that critical thinking matters, but to describe preference aggregation in environments in which a portion of the population is subject to stereotypes. The proofs of the main results, and some immediate extensions, are relegated to the appendix.

4.1 Two Welfare Benchmarks

The relevant unknown is the distribution of reasoned preferences in a large population (continuous), namely the share $p \in [0, 1]$ of individuals who prefer the outcome 1 to the outcome 0 after completing their critical thinking process. Welfare realizes the distance between the social action a and its target p :¹⁷

$$W(a, p) = -(a - p)^2$$

¹⁷Although the space of reasoned preferences — individuals’ resolution of the moral dilemma — is binary, the policy space is continuous. This corresponds to a situation where the planner can fine-tune the policy to the *distribution* of individuals’ reasoned preferences. The example in the introduction of choosing the size of the welfare program based on the share of people who hold an egalitarian (rather than a free market) view fits this story. A different specification would $a^* = \mathbb{I}[p > \frac{1}{2}]$ (binary action space) provide similar insight but is less tractable.

Ex-ante, p is unknown and drawn from a normal distribution $p \sim \mathcal{N}(\mu, \sigma)$; absent the information from the election, the principal would then choose $a = \mu$ and obtain the value $-\sigma^2$.¹⁸ Before choosing $a \in [0, 1]$, the principal observes the proportion \bar{p} of agents that report preferring the alternative 1. We call \bar{p} the *election outcome*. and consider two types of principals that differ in the use they can make of this information.

Positive Principal. A Positive principal, for which the election outcome is *not* binding, namely who can choose any $a \in [0, 1]$ regardless of the implementation of \bar{p} . The positive principal uses the election result and his knowledge of the critical thinking process within the population to estimate p . His optimal action is the conditional expectation

$$a^* = \hat{p} := \mathbb{E}[p|\bar{p}]$$

that achieves value;

$$W_p = -\mathbb{E} \left[(\hat{p} - p)^2 \right], \quad (1)$$

equal to the dispersion of the conditional mean \hat{p} around p . Both expectation operators \mathbb{E} integrate under the joint distribution of p, \bar{p}, \hat{p} , which depend on the voting behavior and citizens' critical thinking process — which we derive in the next section. Connecting to the discussion in the introduction, one can think of such principals as public figures (e.g., multinational firms or social influencers with reputational concerns) who need to take a stance on a dilemma. They privately run a poll and use its outcome as they wish to fine-tune their statement. Payoff depends on the (distribution of) reasoned preferences because the statement acts as a “focusing event” that pushes the relevant population into critical thinking: the preferences individuals judge the principal on are (potentially) different from those they report at the poll.

Institutional Principal. Second, we consider an Institutional principal who has to choose $a = \bar{p}$. One can think of such principals as democratic institutions that must comply with the election outcome (say, by empowering a parliament whose composition is proportional to \bar{p}).¹⁹ Due to the constraint in her action, the Institutional princi-

¹⁸The normality assumption gives tractable conditional expectations and closed-form welfare. It is inconsistent with the compact support $[0, 1]$. The analysis with ex-ante uniform p (and p_S) is algebraically more involved, but does not change the qualitative results. For tractability, we keep the normal setup, implicitly assuming that σ is “small enough” that the mass outside $[0, 1]$ is negligible.

¹⁹In this context, the interpretation of p differs. Rather than focusing on the potential backlash from reasoned preferences, we envision an institutional principle considering p as a normative criterion for aggregating social preferences about a dilemma. Essentially, the distribution of preferences of individu-

pal achieves value.

$$W_I = -\mathbb{E} \left[(\bar{p} - p)^2 \right] \quad (2)$$

via the standard decomposition we obtain

$$W_I = W_P - B, \quad (3)$$

where

$$B = \mathbb{E} \left[(\bar{p} - \hat{p})^2 \right] > 0$$

is the bias of election, representing how the average reported preference differs systematically from the reasoned ones. A principal P who can correct for such social tendencies only suffers from the dispersion of the estimator \hat{p} around the parameter p , while the principal I must also be concerned with the bias of the election.

4.2 Cognitive and Voting Processes

Each individual is characterized by a reasoned preference

$$y \sim \text{Ber}(p)$$

where p is the unknown welfare relevant to which the principal wants to match. For example, an individual with $y = 1$ has a reasoned preference for the alternative 1. However, if asked at a poll, individuals do not necessarily report their reasoned preference. This is because the reasoned preference is “discovered” at the end of a *critical thinking process* that individuals undergo.

The Cognitive Process. Agents transition through two critical thinking states $\{S, A\}$, where S means *Stereotype* and A means *Awareness*. We assume that the critical thinking process follows a simple dynamic in continuous time: all individuals start at $t = 0$ in state S and, independently of y (and other voting parameters), transition to the absorbing state A with intensity $\lambda \in (0, \infty)$. Therefore, at time t there will be a fraction

$$\eta_S = \exp \{-\lambda t\}$$

als who have undergone the critical thinking process determines the “right thing to do”.

of agents that are still Stereotypes and $\eta_A = 1 - \eta_S$ that transitioned to Awareness.²⁰ The parameter λ is key for our analysis. It represents the intensity with which individuals realize that the issue at hand is ambivalent. In our experiment, we established the way news are presented (media) has an effect on λ and that this effect depends on the cognitive abilities of the population. Notice that In models where the principal cares exclusively about the share of critical thinkers, then the result is trivial as η_A is increasing in λ for all t . We show how

Voting Behavior. We denote x the preference that individuals report in the polls and assume it depends on the reasoned preference y and on the stage of the critical thinking process $\{S, A\}$. Before realizing that the issue is ambivalent, the preference reported x_S is

$$x_S | y = \begin{cases} \text{Ber}(p_S) & \text{w.p. } 1 - \beta \\ y & \text{w.p. } \beta \end{cases}$$

In other words, x_S is equal to the reasoned preference with probability $\beta \in [0, 1]$, while the complementary probability is drawn from a distribution of stereotypical preferences $p_S \sim \mathcal{N}(\mu, \sigma)$, independent of p . Since we still have a parameter β driving the correlation between average stereotypes and reasoned preferences, the assumption of independence is innocuous. It only requires the formation of stereotypical preferences involving factors not solely related to p .²¹ Notice that a high β represent situations where, despite not realizing the ambivalent nature of the issue, stereotypes get their reasoned preference right with high probability.²² On the contrary, if $\beta = 0$, corresponding to a situation where stereotypes are independent of reasoned preferences, then an election held at $t = 0$ (all stereotypes) would result in $\bar{p} = p_S$, hence it is not informative at all about p .

²⁰The assumption that A is an absorbing state, with no transitions from A to S , captures the idea that awareness is an irreversible process. A straightforward extension of the model prevents a scenario where all individuals eventually reach state A : a constant fraction $\nu < \lambda$ exits the economy and reenters in the awareness state S . Qualitative results would remain unchanged as the associated share of stereotypes: $\eta_S(t) = \frac{\nu}{\lambda} + \exp(-\lambda t) (1 - \frac{\nu}{\lambda})$ would still be decreasing in λ, t .

²¹The identical distribution of p, p_S is instead for tractability alone. Most derivations in the Appendix utilize nonidentically distributed normal variables $(\mu, \sigma, \mu_S, \sigma_S)$. We discuss such extensions, focusing on the meaning of $\mu \neq \mu_S$, in Section 4.3.

²²Recall that in our setting there is no intrinsic social value for being critical thinkers so if all agents get their reasoned preference right we have perfect elections. However, a related phenomenon studied by [Bernheim et al. \(2021\)](#), “mental flexibility” might have social benefits beyond increasing the accuracy of elections. The challenge for us is to derive λ as a welfare measure even without a direct beneficial effect of critical thinking.

The preference reported by individuals in A loses its dependence on the nuisance parameter p_S and becomes a function of the reasoned preference alone,

$$x_A | y = \begin{cases} y & \text{w.p. } \xi \\ 1 - y & \text{w.p. } 1 - \xi \end{cases}$$

The parameter $\xi \in [\frac{1}{2}, 1]$ is meant to capture situations in which citizens realize that the issue is ambivalent but have not yet found their reasoned preference. The case $\xi = 1$ corresponds to a situation where individuals discover their reasoned preference immediately after realizing the ambivalence of the issue, while $\xi = \frac{1}{2}$ is a situation of permanent indecisiveness of A individuals. We think indeed of our two-stage critical thinking process as a reduced form of a fully identified three-stage process – detailed in the appendix – where A is an intermediate stage where agents have realized the ambivalent nature of the issue but have not formed their reasoned preference yet, i.e., they are in a phase of normative uncertainty. Note that if all individuals are in the state A (that is, a poll held at $t \rightarrow \infty$), then the election result is $\bar{p} = \xi \cdot p + (1 - \xi) \cdot (1 - p)$, which is a strictly monotonic (hence invertible) function of p if $\xi > \frac{1}{2}$. In that case, $\hat{p} = p$, and the *Positive* principal chooses efficiently.²³ For interior shares η , the election outcome is given by:

$$\bar{p} = \eta_S (\beta p_S + (1 - \beta) p) + \eta_A (\xi \cdot p + (1 - \xi) \cdot (1 - p)) \quad (4)$$

and the parameter ξ affects the *Positive* welfare too. We can now use the (joint) normality assumption to write \bar{p} and the conditional expectation \hat{p} as a linear function of the fundamental unknowns p, p_S , that is,

$$\bar{p} = \alpha_0 + \alpha_1 \cdot p + \alpha_2 \cdot p_S$$

$$\hat{p} = \gamma_0 + \gamma_1 \cdot p + \gamma_2 \cdot p_S$$

where loadings α, γ are functions of the structural parameters $\vartheta = [\beta, \xi, \mu, \sigma]$ and the statistic of the critical thinking process η (see appendix). Once we specify the joint normal expectation operator, we can compute (the evolution of) both *Positive* and *Institutional* welfare 2-3 in closed form and arrive at our main result.

Proposition 1 *i) For all values of structural parameters ϑ , W_P is increasing in t and λ .*

ii) W^I has nontrivial comparative statics in λ, t . If $\beta < 1 - \xi$, then it is monotonically

²³Institutional principal still needs to consider the attenuation bias driven by A 's indecisiveness.

increasing; if $\beta > \frac{(1-\zeta)((1-2\mu)^2+4\sigma^2)}{2\sigma^2}$ then it is monotonically decreasing; else it grows locally to $t = 0$ (resp. $\lambda = 0$) up to a finite time t^* (finite intensity λ^*) then eventually decreases.

Figure [to be added] gives a graphical representation of the results collected in Proposition 1, which we now discuss. The point *i* establishes that if the principal knows the value of the structural parameters ϑ and can utilize the outcome of elections without constraints, then the faster individuals move into critical thinking (the higher λ), the higher the efficiency of the elections. Indeed in all the plots of the figure, we observe that positive welfare W^P is increasing in time.²⁴ The reason behind this result is simple to grasp: as fewer and fewer individuals are S , the election outcome \bar{p} becomes less and less dependent on the nuisance unknown p_S which confounds the inference of the welfare relevant unknown p .²⁵ This is an important result for our analysis as it establishes that λ is a welfare measure in a well-definite sense in our setting.

However, in point *ii*), we also hint at a potential limitation of such a result in the case where the principal is constrained to act according to the election outcome due to the (potentially perverse) effect that the movement into critical thinking has on the bias of the election. The most paradoxical result – the condition that if β is large enough, institutional welfare actually *decreases* in λ – has a natural explanation. When β is large, then stereotypes are strong predictors of reasoned preference (at the extreme where $\beta = 1$, all stereotypes vote y despite not realizing the ambivalent nature of the issue), hence moving in the Awareness state indecisiveness and associated attenuation bias, case $\zeta < 1$ – pushes \bar{p} away from p and thus reduces efficiency. This seems – at least to us – a pathological case since it requires. but is useful to highlight the potential role of the bias. For this reason, we further investigate conditions under which the two rules coincide, that is, whether there is a level of η such that “by divine coincidence” the loadings $\alpha = \gamma$ so that the positive and institutional principal have the same action rule – and hence the same value at potential limitations of this interpretation; we indeed show that if $\zeta = 1$ or $\beta < \frac{1}{2}$, the condition for W^I monotonically increasing is vacuously satisfied.

Solving the system of equations $\alpha \equiv \gamma$ gives a share of stereotypes η^* such that the two coincide. Therefore, there exists an interior time where the average reported preference is unbiased for p . Formalizing this result we obtain:

²⁴As the proof relies on W^P being decreasing in the share of stereotypes η , the same graph would be obtained if we fix the time and let λ vary. The bottom-right panel explains the dynamics for high and low λ .

²⁵Indeed, this result does not require the normality assumption but can be directly deduced by the expression of \bar{p} .

Proposition 2 *If there is no bias in the stereotype pool and β is large enough, i.e. if*

$$\mu_p = \mu_{ps} \text{ and } \frac{1 - \beta}{\beta} < \frac{\sigma_{ps}^2}{\sigma_p^2}$$

, then there exists a finite time t^ such that $B(t^*) = 0$. If, in addition,*

$$\xi = 1 \text{ then } t^* = -\frac{1}{\lambda} \log \left(\frac{\sigma_x^2}{\beta (\sigma_x^2 + \sigma_y^2)} \right)$$

with immediate comparative statics.

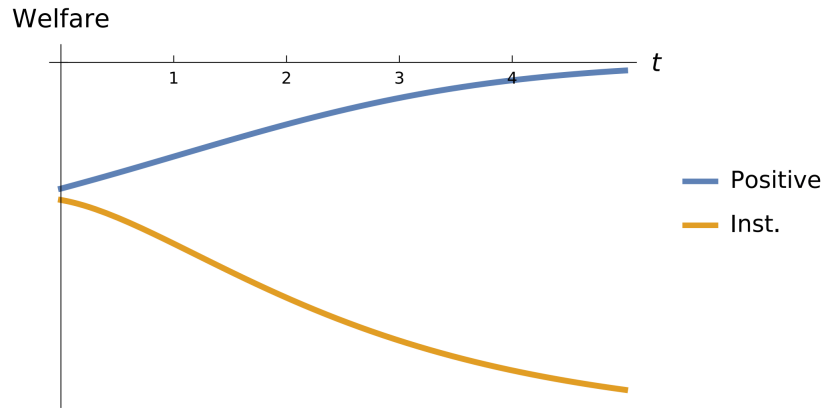


Figure 2: $\beta < 1 - \xi_{NU} \Rightarrow \eta^* = 0 \Rightarrow$ Inst. Welfare is decreasing

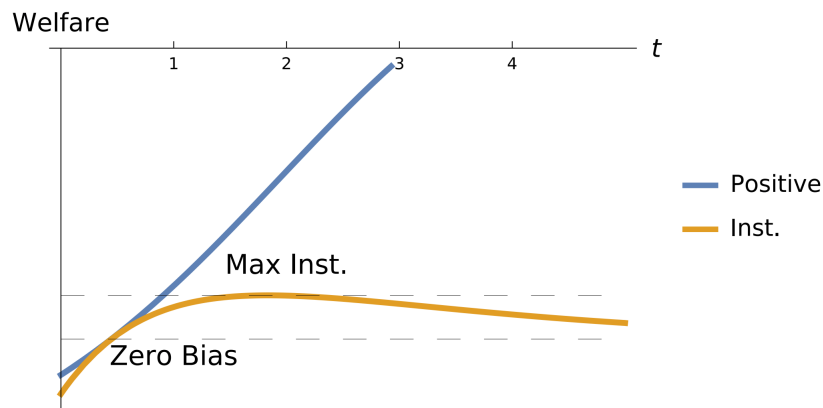


Figure 3: $\eta^* \in (0, 1) \Rightarrow$ Inst. Welfare has interior maximum, after the zero-bias time t^*

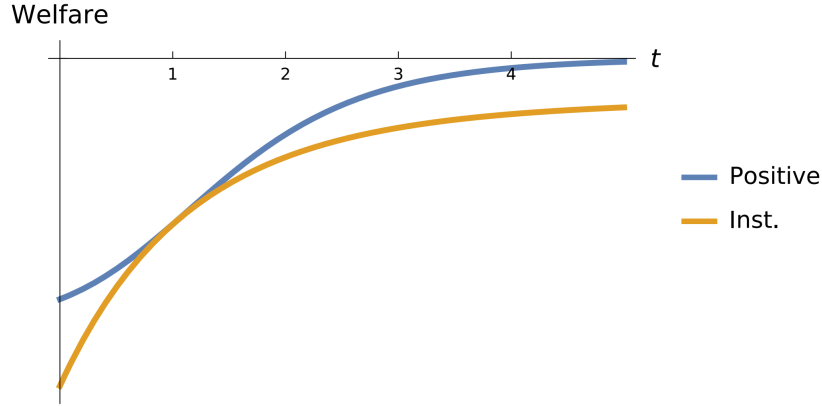


Figure 4: $\eta^* = 1 \Rightarrow$ Inst. Welfare is increasing

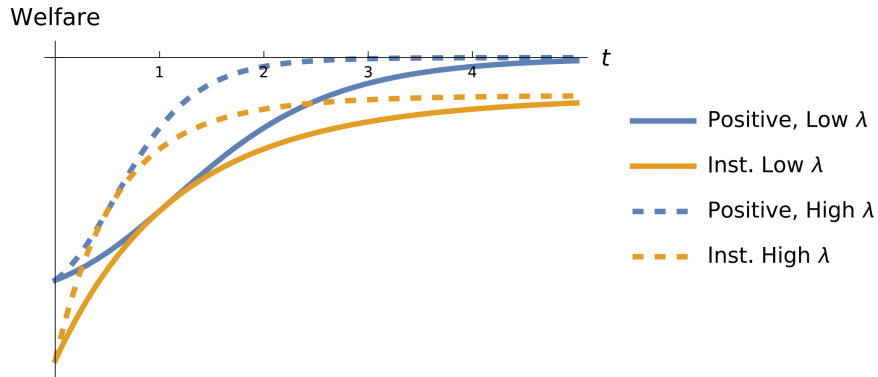


Figure 5: λ is a welfare measure.

4.3 Discussion and Extensions

We have presented a relatively parsimonious model of voting while undergoing a critical thinking process (from Stereotypes to *Aware* citizens). Its mechanics are relatively straightforward: As more citizens become critical thinkers the election outcome is less influenced by the nuisance parameter p_s so the principal can estimate the relevant parameter p more efficiently. Since the share of critical thinkers increases (in every period) with λ , a faster critical thinking process is typically beneficial for the efficiency of elections. We now discuss possible weakening of our model’s assumptions to address potential asymmetries and indicate some adjustment margins that, in light of empirical evidence, one should consider in a richer theoretical model of preference aggregation in the presence of a portion of the electorate who has not even realized the ambivalent nature of the problem.

Three-states cognitive model

Asymmetric settings We have always maintained an implicit assumption of symmetry: the voting structural parameters β, ζ are independent of the reasoned (or stereotypical) preference. A relaxation of this assumption would require modeling $\beta_i = \mathbb{P}[x_S = y|y = i], \zeta_i = \mathbb{P}[x_A = y|y = i]$ with a different specification for the residual uncertainty in the preference of stereotypes.²⁶ Insofar as overconfidence can be interpreted as individuals’ resistance to critical thinking, evidence in [Ortoleva and Snowberg \(2015\)](#) also questions the fact that intensity λ is independent of y : if the reasoned preference predicts cognitive traits associated with critical thinking (or the impact of different storytelling formats), then the *Aware* pool would be selected based on y , which constitutes an additional source of bias.

A similar extension of our model is to allow the presence of bias in the stereotype pool, i.e., to let $\mu_s \neq \mu$, corresponding to a situation where the principal knows that a specific opinion is prevalent before individuals realize the ambivalent nature of an issue. This possibility — which seems compelling whenever one of the positions is more prone to be defended by means of superficial arguments (nationalism) — means the Institutional principal additionally benefits from increasing the intensity λ (or simply “letting time pass”) as having a larger share of *A* voters would mechanically remove this type of systematic bias.²⁷ Extending the model to allow for either of these asymmetries would not alter our main message: If there are more *Aware* voters, polls contain more information about the distribution of reasoned preferences. This is all that matters for a principal who can “filter out” all systematic tendencies in voting, including the asymmetries in stereotype reporting and critical thinking transition, while a principal who cares about getting the election outcome as close as possible to p needs to trade off accuracy with election bias. However, those asymmetries might play a crucial role in determining the voting behavior of critical thinkers who have not discovered their stable preference yet;²⁸ assuming they are aware of such asymmetries – which is reasonable, as they just escaped *Stereotype* state – their problem becomes particularly interesting. Because they have “lost” their stereotype preference and have not formed a reasoned one yet, they are in a state of “normative uncertainty”. If they have

²⁶The symmetry hypothesis can readily be tested in experiments like ours where we observe individuals before starting their critical thinking process and after discovering their stable preference – i.e. using is in the three cognitive states extension of the model –. In our specific setting we could not run such test as the number of subjects that were classified as reasoned preferences at the end of the experiment was extremely small ($N=XXX$), and many were critical thinkers even before starting the treatment, hence any statistical test would have no power.

²⁷The evolution of welfare for the *Positive* principal would instead be unaffected by this extension, as she could “clear out” all systematic noise in the poll.

²⁸Recall our two-state model subsumed this form of indecisiveness in a low value of the ζ parameter.

to vote in such state they might try to compensate the bias in the electorate driven by the asymmetry in the type of voters that get out of the Stereotype phase. That is, they might express preference 0 just because they think that such preference is less prevalent among stereotypes ($\mu_S > \mu$) or because individuals holding that stereotype become critical thinkers faster. Moreover, because those voters don't have a clear preference, they might be especially susceptible to increasing the cost of voting and decide to abstain.²⁹

Discounting Finally, by adding a penalty for waiting (discounting the utility from taking an action later), we can use our setting to discuss the optimal timing of elections. Notice that even the *Positive* principal – who chooses efficiently in the limit – would not postpone her decision until $t = \infty$ under discounting. Studying how the timing of the optimal election varies with intensity λ (and other structural parameters) amounts to analyzing the problem of a principal who controls the type and duration of storytelling to which she wants to subject her agents before administering a poll to maximize its accuracy. Such problems arise naturally in many realistic settings, e.g. designing the type and accuracy of information to give a focus group before asking their opinion on some marketing campaign or policy proposal.

5 Conclusion

We experimented and determined that the format in which the news is presented affects the transition of a person to *A*. This effect is driven by individuals with a high need for cognition (the flexibility scale is insignificant). Realizing the ambivalent nature of an issue is an essential step in discovering a reasoned preference, as it improves the “quality” of one’s preference from raw to reasoned. As such, critical thinking is *also* good for the efficiency of elections.

What is broadly referred to as a storytelling format (e.g., newspapers, television, social media, social echo chambers) might impact the probability of realizing ambivalence. Beyond “informing” and “persuading,” it also affects an individual’s critical thinking process. Additionally, the format in which news is presented—many short messages vs. more coherent but greedy attention discourse—matters. In particular, unexplored physiological drivers were correlated with standard metrics of cog-

²⁹The effects of voting costs that are heterogeneous based on demographics has been shown to be important (Cantoni and Pons??, others!!); we highlight a potential alternative channel through its heterogeneous impact on individuals who are at different phases of their critical thinking process.

nitition/flexibility.

Reasoned preferences y are *not observable*, and the model is not (fully) identified. We rely on a reduced form for an identified model with three cognitive stages $S \rightarrow A \rightarrow T$ and a final transition to a reasoned preference (resolving awareness) with qualitatively similar results. First, $= A$ realized that the issue was ambivalent, but still did not find our y . Second, how did the voters in A vote? Strategic voting in the presence of stereotype bias, “I still have not resolved my awareness about [topic], but I see a lot of prejudice in favor of position 0, so I vote 1 to compensate.” Additionally, the reasoned preference y is independent of other individual types $(\beta, \xi_A, \lambda \dots)$. Prejudices often coincide with a reasoned preference if the latter is 1. $\beta_1 > \beta_0$. Prejudices might be correlated with the likelihood of becoming aware (Ortoleva and Snowberg (2015)).

References

- Allcott, Hunt, Luca Braghieri, Sarah Eichmeyer, and Matthew Gentzkow (2020), “The welfare effects of social media.” *American Economic Review*, 110, 629–676.
- Bénabou, Roland and Jean Tirole (2006), “Incentives and prosocial behavior.” *American Economic Review*, 96, 1652–1678.
- Bernheim, B Douglas, Luca Braghieri, Alejandro Martínez-Marquina, and David Zuckerman (2021), “A theory of chosen preferences.” *American Economic Review*, 111, 720–54.
- Cacioppo, John T and Richard E Petty (1982), “The need for cognition.” *Journal of personality and social psychology*, 42, 116.
- Eliaz, Kfir and Ran Spiegler (2020), “A model of competing narratives.” *American Economic Review*, 110, 3786–3816.
- Falck, Oliver, Robert Gold, and Stephan Heblich (2014), “E-lections: Voting behavior and the internet.” *American Economic Review*, 104, 2238–2265.
- Feddersen, Timothy and Wolfgang Pesendorfer (1997), “Voting behavior and information aggregation in elections with private information.” *Econometrica: Journal of the Econometric Society*, 1029–1058.
- Gentzkow, Matthew and M Jesse Shapiro (2010), “Ideological segregation online and offline.” *National Bureau of Economic Research*, Working Paper 15916.

- Gorodnichenko, Yuriy, Tho Pham, and Oleksander Talavera (2021), "Social media, sentiment and public opinions: Evidence from #brexit and #uselection." *European Economic Review*, 136.
- Gul, Faruk and Wolfgang Pesendorfer (2009), "Partisan politics and election failure with ignorant voters." *Journal of Economic Theory*, 144, 146–174.
- Halpern, Diane F (2013), *Thought and knowledge: An introduction to critical thinking*. Psychology Press.
- Kaplan, Kalman J (1972), "On the ambivalence-indifference problem in attitude theory and measurement: A suggested modification of the semantic differential technique." *Psychological Bulletin*, 77, 361–372.
- Kim, Jaehoon and Mark Fey (2007), "The swing voter's curse with adversarial preferences." *Journal of Economic Theory*, 135, 236–252.
- Kunda, Ziva (1990), "The case for motivated reasoning." *Psychological bulletin*, 108, 480.
- List, A John (2022), "Enhancing critical thinking skill formation: Getting fast thinkers to slow down." *The Journal of Economic Education*, 53, 100–108.
- Martin, Matthew M and Rebecca B Rubin (1995), "A new measure of cognitive flexibility." *Psychological reports*, 76, 623–626.
- Munir, Saba (2018), "Social media and shaping voting behavior of youth: The scottish referendum 2014 case." *The Journal of Social Media in Society*, 7, 253–279.
- Ortoleva, Pietro and Erik Snowberg (2015), "Overconfidence in political behavior." *American Economic Review*, 105, 504–35.
- Shiller, Robert J (2017), "Narrative economics." *American economic review*, 107, 967–1004.
- Vogels, Emily A and Monica Anderson (2019), "Americans and digital knowledge."

Appendices

A Proofs of The Main Model	29
A.1 Preliminary Results on \hat{p} and alike	29
A.2 Proof of Proposition 1	33
A.3 Proof of Proposition 2	35
B Experimental Design Details	37
B.1 Treatments Details and Examples	37
B.2 Detailed Data Collection	38
B.3 Detailed Elicitations	39
B.4 Detailed Description of Graders' Instructions	41
B.5 Heterogeneous Critical Thinking Classification	42
B.6 Threshold changes	45
B.7 AI Digital Grade	47
C Model With Three-State Critical Thinking	47
C.1 Model Identification	49
C.2 General Results	49
C.3 Proofs	57
D Experiment With A Three Cognitive-State Model	57

A Proofs of The Main Model

A.1 Preliminary Results on \hat{p} and alike

The following three steps explicitly show how to analyse the evolution of the cognitive state process over time for each agent and how this relates to the parameters of the model.

1) $\mu_S = \exp\{-\lambda_1 t\}$ and $\mu_C = 1 - \mu_S$ represent the masses. λ_1 represents the intensity with which agents pass from the cognitive state S to the cognitive state A over time. Moreover, we define the unknown parameter \bar{p} as function of μ and p

$$\begin{aligned}
\bar{p}(\mu, p) &= \mu_S (\mathbb{E}[x_S | p]) + \mu_C (\mathbb{E}[x_C | p]) \\
&= \mu_S (\beta p_S + (1 - \beta) p) + \mu_C (\xi_C p + (1 - \xi_C) (1 - p)) \\
&= \mu_S (\beta p_S + (1 - \beta) p) + \mu_C (1 - p - \xi_C (1 - 2p))
\end{aligned}$$

Thus

$$\bar{p}(\mu, p) = \mu_S (\beta p_S + (1 - \beta) p) + \mu_C (1 - p - \xi_C (1 - 2p)) \quad (5)$$

From which we can derive the expression for p as a function of p_S

$$\begin{aligned}
\bar{p} &= \mu_S \beta p_S + \mu_S (1 - \beta) p + \mu_C - \mu_C p - \mu_C \xi_C + 2\mu_C p \xi_C \\
\bar{p} &= \mu_S \beta p_S + \mu_C - \mu_C \xi_C + p [\mu_S (1 - \beta) - \mu_C (1 - 2\xi_C)] \\
&= \mu_S \beta p_S + \mu_C - \mu_C \xi_C + p [\mu_S (1 - \beta) - \mu_C (1 - 2\xi_C)]
\end{aligned}$$

Thus p is defined as

$$p = \frac{\bar{p} - \mu_S \beta p_S - \mu_C (1 - \xi_C)}{\mu_S (1 - \beta) - \mu_C (1 - 2\xi_C)} \quad (6)$$

Finally, we can define the parameter \hat{p} that is defined as the expectation of p conditioning on \bar{p}

$$\hat{p} = \frac{\bar{p} - [\mu_S \beta \mathbb{E}[p_S | \bar{p}] + \mu_C (1 - \xi_C)]}{\mu_S (1 - \beta) + \mu_C (2\xi_C - 1)} \quad (7)$$

Thus \bar{p} is defined as

$$\hat{p} = \mathbb{E}[p | \alpha_1 p + \alpha_2 p_S = \bar{p}] \quad (8)$$

$$\begin{aligned}
\hat{p} &= \frac{\beta^2 \mu_S^2 \mu_X \sigma_Y^2 + \sigma_X^2 (1 - \bar{p} - \xi_C + \mu_S (-1 + \beta \mu_Y + \xi_C)) (1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C))}{\beta^2 \mu_S^2 \sigma_Y^2 + \sigma_X^2 (1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C))^2} \\
&= \frac{\beta^2 \mu_S^2 \mu_X \sigma_Y^2 + \sigma_X^2 (\bar{p} - [\mu_S \beta \mu_Y + (1 - \mu_S) (1 - \xi_C)]) (1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C))}{\beta^2 \mu_S^2 \sigma_Y^2 + \sigma_X^2 (1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C))^2}
\end{aligned}$$

2) It is worth noting that the NWF can be expressed as the sum of the PWF and a biased term due to the elections. Indeed,

$$\begin{aligned}
NWF &= -\mathbb{E} \left[(p - \bar{p})^2 \right] = -\mathbb{E} \left[(p - \hat{p} + \hat{p} - \bar{p})^2 \right] \\
&= -\mathbb{E} \left[(p - \hat{p})^2 + 2(p - \hat{p})(\hat{p} - \bar{p}) + (\hat{p} - \bar{p})^2 \right] \\
&= - \left[\mathbb{E} \left[(p - \hat{p})^2 \right] + 2\mathbb{E} \left[(p - \hat{p})(\hat{p} - \bar{p}) \right] + \mathbb{E} \left[(\hat{p} - \bar{p})^2 \right] \right] \\
&= - \left[\mathbb{E} \left[(p - \hat{p})^2 \right] + 2(\hat{p} - \bar{p}) \underbrace{\mathbb{E} \left[(p - \hat{p}) \right]}_0 + \mathbb{E} \left[(\hat{p} - \bar{p})^2 \right] \right] \\
&= - \left[\underbrace{\mathbb{E} \left[(p - \hat{p})^2 \right]}_{\text{Precision of elections}} + \underbrace{\mathbb{E} \left[(\hat{p} - \bar{p})^2 \right]}_{\text{Bias of elections}} \right]
\end{aligned}$$

Thus it can be rewritten as

$$NWF = PWF + Bias$$

At this stage, we define the two welfare functions given the distributions of the parameters

3) What the theoretical analysis wants to show is the evolution of the welfare functions over time and the main differences between the evolution of the PWF and the NWF. In particular, in order to study the evolution, we take the first derivative of the two functions with respect to μ_S . It is necessary and sufficient to show the sign of this derivative in order to have an all-rounded understanding of the evolution of the two functions. In fact, μ_S as defined above depends negatively on t and λ_1 . Hence, once we define the relation between the functions and μ_S , we immediately get to know the relation between the functions and the time/lambda. Thus, let us start by showing the behavior of the PWF.

$$\frac{\partial PWF}{\partial \mu_s} = \frac{2\beta^2 \mu_s \sigma^2 (-1 + 2\zeta_C) [1 - 2\zeta_C + \mu_s (-2 + \beta + 2\zeta_C)]}{\left\{ 2\mu_s^2 \left[\beta^2 + 2\beta(-1 + \zeta_C) + 2(-1 + \zeta_C)^2 \right] + (1 - 2\zeta_C)^2 - 2\mu_s(-1 + 2\zeta_C)(-2 + \beta + 2\zeta_C) \right\}^2} \propto 1 - 2\zeta_C$$

Since almost everything is bigger or equal than 0, if we want to study the sign of the above formula, then we just have to analyse the sign of the following term

$$1 - 2\zeta_C + \mu_s (-2 + \beta + 2\zeta_C) < 0$$

$$0 < \mu_s < \underbrace{\frac{2\zeta_C - 1}{-2 + \beta + 2\zeta_C}}_{\geq 1?}$$

Proposition 3 W_P is increasing in t and λ if

$$\frac{2\zeta_C - 1}{-2 + \beta + 2\zeta_C} > 1 \iff 1 > \beta$$

Let's study the right hand side of the inequality

$$2\zeta_C - 1 \geq -2 + \beta + 2\zeta_C$$

$$\beta \leq 1$$

Therefore, we can conclude that PWF is decreasing in μ_s for each time t , because

$$0 < \mu_s < \frac{2\zeta_C - 1}{-2 + \beta + 2\zeta_C}, \quad \forall \mu_s \in [0, 1]$$

In other words, the PWF is an increasing function of both t and λ_1 .

A.2 Proof of Proposition 1

Proof 1 Notice preliminary that using the chain rule the following result is valid for both welfare functions

$$\frac{dW}{d\lambda} = \frac{dW}{d\eta} \cdot \underbrace{\frac{d\eta}{d\lambda}}_{<0} \implies \frac{dW}{d\lambda} \propto -\frac{dW}{d\eta}$$

and therefore welfare moves in λ (and t) contrary to how it moves in the share of stereotypes. First, for positive welfare, $\frac{dW}{d\eta}$ is always negative for the following computations

$$\begin{aligned} \frac{\partial PWF}{\partial \mu_S} = & \frac{2\beta^2 \mu_S \sigma^2 (-1 + 2\xi_C) [1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C)]}{\left\{ 2\mu_S^2 \left[\beta^2 + 2\beta(-1 + \xi_C) + 2(-1 + \xi_C)^2 \right] + (1 - 2\xi_C)^2 - 2\mu_S(-1 + 2\xi_C)(-2 + \beta + 2\xi_C) \right\}^2} \\ & \propto 1 - 2\xi_C + \mu_S (\beta - 2(1 - \xi_C)) \end{aligned}$$

Since almost everything is bigger or equal than 0, if we want to study the sign of the above formula, then we just have to analyse the sign of the following term

$$\begin{aligned} 1 - 2\xi_C + \mu_S (-2 + \beta + 2\xi_C) &< 0 \\ 0 < \mu_S &< \underbrace{\frac{2\xi_C - 1}{-2 + \beta + 2\xi_C}}_{\geq 1?} \end{aligned}$$

Let's study the right hand side of the inequality

$$\begin{aligned} 2\xi_C - 1 &\geq -2 + \beta + 2\xi_C \\ \beta &\leq 1 \end{aligned}$$

Therefore, we can conclude that PWF is decreasing in μ_S for each time t , because

$$0 < \mu_S < \frac{2\xi_C - 1}{-2 + \beta + 2\xi_C}, \quad \forall \mu_S \in [0, 1]$$

Furthermore, $\frac{dW}{d\eta}$ has a non-trivial solution. That is by studying the sign of the derivative

of welfare elections with respect to μ_S we obtain

$$\frac{\partial NWF}{\partial \mu_S} = - \left[4\beta^2 \mu_S \sigma^2 + 4\beta (-1 + \mu_S) \sigma^2 (-1 + \xi_C) + 4\beta \mu_S \sigma^2 (-1 + \xi_C) + 2(-1 + \mu_S) \left[(1 - 2\mu)^2 + 4\sigma^2 \right] (-1 + \xi_C)^2 \right]$$

The sign of the term in brackets is

$$4\beta^2 \mu_S \sigma^2 + 4\beta \sigma^2 (-1 + \xi_C) (-1 + 2\mu_S) + 2(-1 + \mu_S) \left[(1 - 2\mu)^2 + 4\sigma^2 \right] (-1 + \xi_C)^2 > 0$$

$$\mu_S \left[4\beta^2 \sigma^2 + 8\beta \sigma^2 (-1 + \xi_C) + 2(-1 + \xi_C)^2 \left[(1 - 2\mu)^2 + 4\sigma^2 \right] \right] > 4\beta \sigma^2 (-1 + \xi_C) + 2(-1 + \xi_C)^2 \left[(1 - 2\mu)^2 + 4\sigma^2 \right]$$

$$\mu_S > \frac{4\beta \sigma^2 (-1 + \xi_C) + 2(-1 + \xi_C)^2 \left[(1 - 2\mu)^2 + 4\sigma^2 \right]}{\underbrace{4\beta^2 \sigma^2 + 8\beta \sigma^2 (-1 + \xi_C) + 2(-1 + \xi_C)^2 \left[(1 - 2\mu)^2 + 4\sigma^2 \right]}_{\text{Threshold} < 1?}}$$

Saying that the threshold is less than one also means that $\frac{\partial NWF}{\partial \mu_S} < 0 \iff \mu_S > \text{Threshold}$. We want to study this threshold. The conditions given in the text correspond to this threshold being below 0 and above 1, respectively.

$$\begin{cases} \text{COND1} \rightarrow \text{Threshold} < 0 & \text{Always decreases in time} \\ \text{COND2} \rightarrow \text{Threshold} > 1 & \text{Always increase in time} \\ \text{COND3} \rightarrow \text{Threshold} \in (0, 1) & \text{Increases first, decreases later} \end{cases}$$

COND1 occurs according to the following expression

$$\beta > \frac{(1 - \xi_C) \left((1 - 2\mu)^2 + 4\sigma^2 \right)}{2\sigma^2}$$

Then the welfare is always decreasing. For COND2 to occur, the numerator of the threshold must be higher than the denominator. Hence, since there are only two terms differing between numerator and denominator, the following must be true.

$$4\beta \sigma^2 (-1 + \xi_C) > 4\beta^2 \sigma^2 + 8\beta \sigma^2 (-1 + \xi_C)$$

$$4\beta \sigma^2 (-1 + \xi_C) > 4\beta \sigma^2 (\beta + 2\xi_C - 2)$$

$$\beta < 1 - \xi_C$$

The welfare is then always increasing. Finally, COND3 can be intuitively discussed. Since μ_S is monotonically decreasing in time, there must be continuity of t^{\max} such that

$$\begin{cases} \frac{\partial W^N}{\partial \mu_S} < 0 & \text{for } t < t^{max} \\ \frac{\partial W^N}{\partial \mu_S} < 0 & \text{for } t > t^{max} \end{cases}$$

Therefore, t^{max} is a maximum interior of W^N when threshold $\in (0, 1)$

A.3 Proof of Proposition 2

Proof 2 Firstly, define the parameters associated with \bar{p}

$$\bar{p} = \mu_T p + \mu_S (\beta p_S + (1 - \beta) p) + \mu_C [1 - p + \xi_C (2p - 1)]$$

where

$$\begin{aligned} \alpha_0 &= \mu_C (1 - \xi_C) \\ \alpha_1 &= 1 - \beta \mu_S - 2\mu_C (1 - \xi_C) \\ \alpha_2 &= \beta \mu_S \end{aligned}$$

Then \hat{p} is given by

$$\hat{p} = \frac{\frac{\bar{p} - [\alpha_0 + \alpha_2 \mu_y]}{\alpha_1} \alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

where

$$\begin{aligned} \gamma_0(t, \lambda) &= \frac{\alpha_2^2 \sigma_y^2 \mu_x - \alpha_2 \alpha_1 \sigma_x^2 \mu_y}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} \\ \gamma_1(t, \lambda) &= \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} \\ \gamma_2(t, \lambda) &= \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} \end{aligned}$$

The bias is zero if and only if the following system has a solution

$$\begin{cases} \alpha_1 = \gamma_1 \\ \alpha_2 = \gamma_2 \end{cases}$$

that is

$$\begin{cases} \alpha_1 = \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} \\ \alpha_2 = \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} \end{cases}$$

It is immediate to check that $\gamma_1(t, \lambda) = \alpha_1(t, \lambda) \iff \gamma_2(t, \lambda) = \alpha_2(t, \lambda)$, so we actually have a single equation and we need to claim that exists a time such that

$$\begin{aligned} \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \alpha_1 &\iff \frac{(1 - \beta \mu_S - 2\mu_C(1 - \xi_C))^2 \sigma_x^2}{(1 - \beta \mu_S - 2\mu_C(1 - \xi_C))^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2} = (1 - \beta \mu_S - 2\mu_C(1 - \xi_C)) \\ &\iff (1 - \beta \mu_S - 2\mu_C(1 - \xi_C)) \sigma_x^2 = (1 - \beta \mu_S - 2\mu_C(1 - \xi_C))^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2 \\ &\iff (1 - \beta \mu_S - 2\mu_C(1 - \xi_C)) \sigma_x^2 (\beta \mu_S + 2\mu_C(1 - \xi_C)) = (\beta \mu_S)^2 \sigma_y^2 \\ &\iff \frac{(1 - \beta \mu_S - 2\mu_C(1 - \xi_C)) (\beta \mu_S + 2\mu_C(1 - \xi_C))}{(\beta \mu_S)^2} = \frac{\sigma_y^2}{\sigma_x^2} \\ &\iff \frac{(1 - \beta \mu_S - 2(1 - \mu_S)(1 - \xi_C)) (\beta \mu_S + 2(1 - \mu_S)(1 - \xi_C))}{(\beta \mu_S)^2} \end{aligned}$$

When $\mu_S = 0$ there cannot be the zero-bias time, because as $t \rightarrow \infty$ this explodes (? can we show this is always increasing in μ_S) because in the limit there is always bias. On the other hand, there could be a zero-bias time that coincides with $t^* = 0$. Indeed, when $\mu_S = 1$

$$\frac{1 - \beta}{\beta} = \frac{\sigma_y^2}{\sigma_x^2}$$

An even more special case is when $\xi_C = 1$. Indeed,

$$\begin{aligned} \frac{(1 - \beta \mu_S)^2 \sigma_x^2}{(1 - \beta \mu_S)^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2} &= (1 - \beta \mu_S) \\ (1 - \beta \mu_S) \sigma_x^2 &= (1 - \beta \mu_S)^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2 \\ \left(\frac{1 - \beta \mu_S}{\beta \mu_S} \right) &= \frac{\sigma_y^2}{\sigma_x^2} \end{aligned}$$

Substituting the expression of μ_S as a function of t and λ

$$\frac{\sigma_x^2}{\beta (\sigma_x^2 + \sigma_y^2)} = e^{-t\lambda_1}$$

that becomes

$$t^* = -\frac{1}{\lambda_1} \log \left(\frac{\sigma_x^2}{\beta (\sigma_x^2 + \sigma_y^2)} \right)$$

where the argument of the log must be smaller than 1

$$\frac{\sigma_x^2}{\beta (\sigma_x^2 + \sigma_y^2)} < 1$$

that is

$$\frac{(1 - \beta)}{\beta} < \frac{\sigma_y^2}{\sigma_x^2}$$

B Experimental Design Details

B.1 Treatments Details and Examples

In the **NEWSPAPER** treatment, the participants are exposed to two news articles: one that is for and one that is against the issue. In the **FACEBOOK** treatment, participants were exposed to six Facebook posts: two for and two against an issue as well as two irrelevant posts. In treatment **TWITTER**, participants are exposed to twenty-four tweets: ten for digital privacy, ten against digital privacy, and four irrelevant tweets. Each tweet has an average length of 40 characters, corresponding to 20 words.³⁰ We give participants 5 seconds to read each tweet before the next one automatically pops until the last one, which is in line with the average reading speed in the US population. In the **PARTISAN TWITTER** treatment, participants are exposed to 13 tweets: 10 for and 3 irrelevant ones or 10 against and 3 irrelevant ones. Within each treatment, tweets, Facebook posts, and news articles arrive in a random order sequentially (one by screen) and remain on screen for a given fixed amount of time (the participant cannot move to the next screen by him or herself). Each participant is randomly assigned to one of the treatments.

see the online appendix about participant's experimental instructions.

³⁰This corresponds to the average length of tweets on twitter.com, see the Appendix for details.

B.2 Detailed Data Collection

Preventing duplicates. Submissions to studies on Prolific are guaranteed to be unique by the firm³¹. Our system is set up such that each participant can have only one submission per study on Prolific. That is, each participant will be listed in your dashboard only once, and can only be paid once. On our side, we also prevent participants to take up several times our experiment in two steps. First, we enable the functionality “Prevent Ballot Box Stuffing” which permits to... Second we check participant ID and delete the second submission from the data set of the same ID if we find any.

Drop-out rates. Here put the drop out (or in the main text).

High vs low-quality submissions. Participants joining the Prolific pool receive a rate based on the quality of their engagement with the studies. If they are rejected from a study then they receive a malus. If they receive too much malus, then they are removed by the pool from the company³². Based on this long term contract, participants are incentivized to pay attention and follow the expectations of each study. Hence, a good research behavior has emerged on Prolific according to which, participants themselves can voluntarily withdraw their submissions if they feel they did a mistake such as rushing too much, letting the survey opened for a long period of time without engaging with it, and so on³³. According to these standards, we kept submissions rejections as low as possible, following standard in online experimental economics. Participants who fail at least one fair attention check are rejected and not paid. Following Prolific standards, participants who are statistical outliers (3 standard deviations below the mean) are excluded from the good complete data set.

Payments and communication. We make sure to review participants’ submissions within within 24-48 hours after they have completed the study. This means that within this time frame, if we accept their submission, they receive their fixed and bonus payment. Otherwise, we reject their submissions and send to them a personalized e-mail⁽³⁴⁾, detailing the reason of the rejection, leaving participants the opportunity to contact us afterwards if they firmly believe the decision to be unfair (motivate their

³¹See Prolific unique submission guarantee policy [here](#).

³²See Prolific pool removal Policy [here](#).

³³See Prolific update regarding this behavior [here](#).

³⁴Partially-anonymized through Prolific messaging app which put the researcher’s name visible to the participants and only the participants ID visible to the researcher.

perspective). Participants can also contact us at any time if they encounter problems with our study or just have questions about it.

B.3 Detailed Elicitations

B.3.1 Political Preferences

Ex-ante and ex-post political preferences. Before the treatment, we prompt participants on different political issues (i.e., without a baseline): guns, crime, climate, welfare, and digital privacy issues. We use the standard congressional metrics, including digital issues. We elicit more than only digital preferences to ensure that participants do not guess at this stage which preferences we focus on in the remaining of the experiment (treatment and critical thinking essay), to minimize their social desirability bias. After the treatment on digital privacy, we survey again participants to elicit their preferences about digital privacy. We use the following scale.

1. On the issue of gun regulation, do you support or oppose the following proposal?
2. On the issue of environmental policies, do you support or oppose the following proposal?
3. On the issue of crime policies, do you support or oppose the following proposal?
4. On the issue of digital policies, do you support or oppose each of the following proposals?

B.3.2 Digital Knowledge Test

see the online appendix about participant's experimental instructions.

B.3.3 Issue Familiarity

1. In the remainder of the experiment, we will focus on the following political issue. Please state again your preference.
2. Have you thought deeply about this issue before participating in this study?
[Yes/No]

B.3.4 Listing Reasons

If *yes* to the previous question, then participants see this question:

You answered “Yes” to the previous question. you will be asked now to provide, at most, two reasons which justify your position and two reasons which justify the opposite position. If you do not know any reasons, please select “I am unable to list any logical reason at the moment”. you do not need to agree with these reasons: they just need to be a logical justification for or against your position. your payment WILL NOT depend on your answer to this question. However, your honest answer is of paramount importance for the success of this study.

1. Reasons which justify your position

- Reason 1: [write text here]
- Reason 2: [write text here]
- I am unable to list any logical reason at the moment

2. Reasons which oppose your position

- Reason 1: [write text here]
- Reason 2: [write text here]
- I am unable to list any logical reason at the moment

B.3.5 Internal Uncertainty

How certain are you of your preference regarding the digital privacy issue? By “Certain”, we mean that you feel confident enough to vote for your political preference if asked to you in a real life political committee. Select among the following options:

- Completely Uncertain
- Rather Uncertainty
- Rather Certain
- Completely Certain

B.3.6 Need for Cognition

For each sentence below, please select how uncharacteristic or characteristic this is for you personally.

B.3.7 Cognitive Flexibility

B.3.8 Habits of News Consumption

see the online appendix about participant’s experimental instructions.

B.4 Detailed Description of Graders’ Instructions

We recruited 20 psychologists (doctoral level or above) who specialize in cognitive psychology at Princeton University. Each grader was randomly assigned a “grading treatment” (that is, a set of essays to grade). Such a set of essays was randomly built, containing essays from all four treatments. Additionally, the graders were not informed about the treatment to which subjects were assigned. Psychologists must grade a very short paragraph (around 300 words or fewer) as follows. The grading consists of giving a passing grade if psychologists judge that the participant “realizes that the issue is ambivalent,” a failing grade otherwise. What may happen is to confound high cognitive sophistication (i.e., the ability to write well-written essays in English), facilitated by the fact that they read some arguments right before this essay exercise with their self-reasoning skill “realizing that the issue is ambivalent”, which is the variable that we want to elicit. This is a specific case that is still challenging for AI-based grading software and the main reason why human expertise is uniquely useful.

We define “realizing that the issue is ambivalent” as the awareness of an individual to recognize that there can be perfectly logical but opposite arguments in favor of and against the same issue that renders the decision-making process complex. Such attitudinal ambivalence leads to temporarily conflicting preferences; namely, one preference for the issue at hand and one preference against the issue at hand. There are different ways of measuring this “awareness,” as documented in the social psychology and cognitive psychology literature. In our study, we capture this awareness by observing individuals reasoning and elaborating in a personal way on the pros and cons of the same issue in a textual format.

Each grader was paid a fixed fee of \$50 for each grading session. Each grader could participate up to three times in our experiment, and no grader could be assigned twice to the same grading treatment. For robustness, each essay was corrected three times by different psychologists. Despite “triple-eliciting” such grades, this metric can still be prone to measurement error. Accordingly, we suggest interpreting the estimated *levels* of λ with caution. However, our focus is on the difference between the treatments. Therefore, such measurement error does not affect this difference.

see the online appendix about the participant’s experimental instructions.

B.5 Heterogeneous Critical Thinking Classification

B.5.1 Critical Thinking Classification Results

Table 4 shows the classification results of individuals as Stereotype and Aware.

Treatment	$S_0 \rightarrow S_1$	$S_0 \rightarrow A_1$	$A_0 \rightarrow A_1$
NEWSPAPER	111	49	12
TWITTER	135	43	15
FACEBOOK	111	60	11
<i>N</i>	357	152	38

Table 4: TABLE 2: CLASSIFICATION RESULTS BEFORE / AFTER TREATMENT

B.5.2 Awareness With Cognitive Styles Heterogeneity with Unanimity

Awareness With Cognitive Flexibility, with Unanimity

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	·	1.301 (0.091)	1.661 (0.095)
TWITTER	·	·	0.422 (0.093)
FACEBOOK	·	·	·

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 5: z-score FOR *High Need for Cognition* WITH UNANIMITY

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	0.320 (0.066)	-0.388 (0.065)
TWITTER	.	.	-0.965 (0.064)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 6: TABLE 6: z-score FOR *Low Need for Cognition* WITH UNANIMITY

Awareness With Need for Cognition, with Unanimity

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	1.061 (0.084)	-2.238* (0.084)
TWITTER	.	.	-1.300 (0.083)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 7: TABLE 6: t-ratio FOR *High Need for Cognition* WITH UNANIMITY

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	0.455 (0.069)	0.705 (0.069)
TWITTER	.	.	0.258 (0.069)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 8: TABLE 6: t-ratio FOR *Low Need for Cognition* WITH UNANIMITY

B.5.3 Awareness With Cognitive Styles Heterogeneity with Majority

Awareness With Cognitive Flexibility, with Majority

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	0.658 (0.076)	-0.924 (0.087)
TWITTER	.	.	-1.609 (0.081)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 9: TABLE 6: *t*-ratio FOR High Flexibility

Treatment	NP	TWITTER	FACEBOOK
NEWSPAPER	.	1.132 (0.062)	-0.388 (0.064)
TWITTER	.	.	-1.564 (0.061)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 10: TABLE 7: *t*-ratio FOR Low Flexibility

Awareness With Need for Cognition, with Majority

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	0.764 (0.070)	-2.238* (0.079)
TWITTER	.	.	-3.087** (0.075)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 11: TABLE 4: *t-ratio* FOR HIGH NEED FOR COGNITION

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.094 (0.066)	0.703 (0.067)
TWITTER	.	.	-0.396 (0.063)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 12: TABLE 5: *t-ratio* FOR LOW NEED FOR COGNITION

B.6 Threshold changes

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.278 (0.054)	-0.923 (0.054)
TWITTER	.	.	-2.262* (0.053)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 13: *t-ratio* DIFFERENCE-IN-MEANS WITH THRESHOLD OF KTS = 8 AND REASON COUNTER = 2

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.222 (0.054)	-0.799 (0.054)
TWITTER	.	.	-2.067* (0.053)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 14: *t-ratio* DIFFERENCE-IN-MEANS WITH THRESHOLD OF KTS = 7 AND REASON COUNTER = 3

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.136 (0.049)	-0.703 (0.052)
TWITTER	.	.	-1.985* (0.049)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 15: *t-ratio* DIFFERENCE-IN-MEANS WITH THRESHOLD OF KTS = 7 AND REASON COUNTER = 2

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.428 (0.049)	-0.507 (0.052)
TWITTER	.	.	-1.985* (0.049)
FACEBOOK	.	.	.

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 16: *t-ratio* DIFFERENCE-IN-MEANS WITH THRESHOLD OF KTS = 6 AND REASON COUNTER = 3

B.7 AI Digital Grade

The last robustness check that we perform is to split the sample conditioning on the the grade of the essay on the digital topic evaluated by the algorithm of the AI. Indeed, the essay is written after undertaking the experiment and it might influence the writing quality of the essay. In particular, our reasoning is that if no difference in proportion is statistically significant, this means that there is no systematic difference between those who were treated through Facebook and those through Twitter, and indeed from Table 17 this is the case.

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	1.716 (0.078)	0.234 (0.081)
TWITTER	.	.	-1.565 (0.074)
FACEBOOK	.	.	.

Standard errors in parentheses
* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 17: *t-ratio* WITH HIGH AI DIGITAL GRADES

Treatment	NEWSPAPER	TWITTER	FACEBOOK
NEWSPAPER	.	0.098 (0.061)	-1.093 (0.067)
TWITTER	.	.	-1.207 (0.065)
FACEBOOK	.	.	.

Standard errors in parentheses
* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 18: *t-ratio* WITH LOW AI DIGITAL GRADES

C Model With Three-State Critical Thinking

We propose an additional model where agents can be at three different states of critical thinking: not engaging with critical thinking, performing critical thinking (either in its first or second state), and having finished performing critical thinking. In our two-stage model, we considered performing critical thinking as having finished performing

it. In this scenario, we propose a three-stage (not fully identified) model that considers the three stages distinctively.

In this economy, the object of interest is the distribution of reasoned preferences over a binary policy space in a large population. Each individual j inside the population is characterized by a three-dimensional type

$$(x_j, y_j, i_j) \in \mathcal{J} := \{0, 1\} \times \{0, 1\} \times \{0, 1\}$$

where x_j , represents the stereotypical preference individual j would self-report when presented with an dilemma for the first time – that is, by definition, before undergoing a critical thinking phase; y_j differs potentially from x_j as it represents the reasoned preference that j holds after completing their period of critical thinking; the cognitive type i_j refers to the cognitive type i_j , interacting with the format, determines how easily individual j moves into (and out of) critical thinking.

Individuals go through a three-step process of “critical thinking” as they form their preferences. The process begins with a “stereotypical-self” state, followed by a period of critical thinking, and ultimately leading to a “reasoned-self” state. We assume that this process is irreversible and that once individuals reach a reasoned-self state, they no longer question their preferences. There is no additional “information” that has to come and change their worldview: the process of critical thinking provides a final and reasoned answer to dilemmas. When asked to report their preferences on a policy issue, individuals in either their stereotypical self or reasoned self state will vote according to their respective preferences, x_j, y_j , respectively. Those who are still in the critical thinking phase will abstain from voting.

The transition between the different phases is determined by an individual’s cognitive style and the characteristics of the storytelling format. Hence, the storytelling format is instrumental in the agent’s transition from a stereotypical state to the reasoned one. By constructing our model, this transition is captured by the critical thinking phase. An economy of reasoned preferences is preferable from efficiency and welfare perspectives to an economy of stereotypical preferences. We formally present such an economy below.

C.1 Model Identification

Using reported preferences of individuals that do the $S(\text{Stereotype}) \rightarrow T(\text{Type})$ transition (i.e. we observe ex ante x_S then y), we get

$$\mathbb{E}[x_S | y = 1] = (1 - \beta) + \beta p_S$$

$$\mathbb{E}[x_S | y = 0] = \beta p_S$$

which gives the estimators

$$\hat{\beta} = 1 - (\bar{x}_{S|1} - \bar{x}_{S|0})$$

and

$$\hat{p}_S = \frac{\bar{x}_{S|0}}{\hat{\beta}}$$

clearly $\hat{p} = \bar{y}$. Finally, using the reported preferences of individuals that do the $A \rightarrow T(\text{Type})$ transition we can estimate ζ_A as

$$\mathbb{E}[x_A | y = 1] = \zeta_A$$

$$\mathbb{E}[x_A | y = 0] = 1 - \zeta_A$$

so $\hat{\zeta}_A = \bar{x}_{NU|1}$ or $\hat{\zeta}_A = 1 - \bar{x}_{NU|0}$. Notice that we can test the assumed symmetry by testing that $\hat{\zeta}_A = \hat{\hat{\zeta}}_A$. Since in our dataset we have few agents that start in A this test has almost no power.

C.2 General Results

The basic decomposition

$$W_E = W_P + \text{Bias}$$

$$-\mathbb{E}[(p - \bar{p})^2] = -\left(\mathbb{E}[(p - \hat{p})^2] + \mathbb{E}[(\hat{p} - \bar{p})^2]\right)$$

is still clearly valid. However, \bar{p} is now given by

$$\begin{aligned} \bar{p} &= \mu_T p + \mu_S (\beta p_S + (1 - \beta) p) + \mu_A [1 - p + \zeta_A (2p - 1)] \\ &= \alpha_0(t, \lambda) + \alpha_1(t, \lambda) p + \alpha_2(t, \lambda) p_S \end{aligned}$$

with

$$\begin{aligned}
\alpha_0 &= \mu_A (1 - \xi_A) \\
\alpha_1 &= 1 - \beta\mu_S - 2\mu_A (1 - \xi_A) \\
\alpha_2 &= \beta\mu_S
\end{aligned}$$

where \hat{p} (that was wrong in the previous file since for non-normal random variables we do not know the expectation of p given the convex combination $\beta p_S + (1 - \beta) p$) is given by

$$\hat{p} = \frac{\frac{\bar{p} - [\alpha_0 + \alpha_2 \mu_y]}{\alpha_1} \alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \tilde{\alpha}_0(t, \lambda) + \tilde{\alpha}_1(t, \lambda) p + \tilde{\alpha}_2(t, \lambda) p_S$$

$$\frac{\frac{\alpha_0(t, \lambda) + \alpha_1(t, \lambda) p + \alpha_2(t, \lambda) p_S - [\alpha_0 + \alpha_2 \mu_y]}{\alpha_1} \alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

so

$$\tilde{\alpha}_1(t, \lambda) = \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\tilde{\alpha}_2(t, \lambda) = \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

when is it

$$\tilde{\alpha}_1(t, \lambda) = \alpha_1(t, \lambda) \iff \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \alpha_1 \iff$$

$$\text{Same } \sigma = \alpha_1 (1 - \alpha_1) = \alpha_2^2 \iff (1 - \beta\mu_S) (\beta\mu_S) = (\beta\mu_S)^2$$

It is immediate to check that $\tilde{\alpha}_1(t, \lambda) = \alpha_1(t, \lambda) \iff \tilde{\alpha}_2(t, \lambda) = \alpha_2(t, \lambda)$ so we actually have a single equation and we need to claim that \exists time such that

$$\begin{aligned} \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \alpha_1 &\iff \frac{(1 - \beta \mu_S - 2\mu_A(1 - \xi_A))^2 \sigma_x^2}{(1 - \beta \mu_S - 2\mu_A(1 - \xi_A))^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2} = (1 - \beta \mu_S - 2\mu_A(1 - \xi_A)) \\ &\iff (1 - \beta \mu_S - 2\mu_A(1 - \xi_A)) \sigma_x^2 = (1 - \beta \mu_S - 2\mu_A(1 - \xi_A))^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2 \\ &\iff (1 - \beta \mu_S - 2\mu_A(1 - \xi_A)) \sigma_x^2 (\beta \mu_S + 2\mu_A(1 - \xi_A)) = (\beta \mu_S)^2 \sigma_y^2 \\ &\iff \frac{(1 - \beta \mu_S - 2\mu_A(1 - \xi_A)) (\beta \mu_S + 2\mu_A(1 - \xi_A))}{(\beta \mu_S)^2} = \frac{\sigma_y^2}{\sigma_x^2} \end{aligned}$$

now substituting μ_S, μ_A we have the LHS is increasing to ∞ in t , therefore there is a unique solution provided that it starts below $\frac{\sigma_y^2}{\sigma_x^2}$, that is if $\frac{1-\beta}{\beta} < \frac{\sigma_y^2}{\sigma_x^2}$ (β is large enough)
as $\mu_S \rightarrow 0$, this explodes (there is always bias in the limit), while at the beginning there is zero bias iff

$$\frac{1 - \beta}{\beta} = \frac{\sigma_y^2}{\sigma_x^2}$$

[example, $\sigma_x^2 = \frac{1}{2}, \sigma_y^2 = \frac{1}{6} \beta = \frac{3}{4} \implies \frac{1-\beta}{\beta} = \frac{1}{3}$]

this is the zero-bias time. Even more special cases $\xi_A = 1$

$$\frac{(1 - \beta \mu_S)^2 \sigma_x^2}{(1 - \beta \mu_S)^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2} = (1 - \beta \mu_S)$$

$$(1 - \beta \mu_S) \sigma_x^2 = (1 - \beta \mu_S)^2 \sigma_x^2 + (\beta \mu_S)^2 \sigma_y^2$$

$$\left(\frac{1 - \beta \mu_S}{\beta \mu_S} \right) = \frac{\sigma_y^2}{\sigma_x^2}$$

Since the LHS is decreasing in μ_S and the RHS is increasing, then there is at most one solution. It has none if

$$\frac{1 - \beta}{\beta} > \frac{\sigma_y^2}{\sigma_x^2}$$

Furthermore we get

$$W_P = -\mathbb{E} \left[(p - \hat{p})^2 \right] = -\frac{\alpha_2^2 \sigma_y^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\text{If } \sigma_x^2 = \sigma_y^2 = \sigma^2 = -\frac{\alpha_2^2}{\alpha_1^2 + \alpha_2^2} \sigma^2$$

$$\xi_A = 1 = \frac{(\beta \mu_S)^2}{2\beta \mu_S [1 - \beta \mu_S] + 1} \sigma^2??$$

Aside: No Bias

Condition for no bias is that coefficients in \bar{p} are the same as in \hat{p} that is,

$$\text{If } \exists t : B(t) = 0$$

$$\alpha_0 + \alpha_1 p + \alpha_2 p_S$$

$$\frac{\frac{\alpha_1 p + \alpha_2 p_S - \alpha_2 \mu_y}{\alpha_1} \alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \frac{\alpha_2^2 \sigma_y^2 \mu_x - \alpha_2 \alpha_1 \sigma_x^2 \mu_y}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} + \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} p + \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} p_S$$

$$\alpha_0 = \frac{\alpha_2^2 \sigma_y^2 \mu_x - \alpha_2 \alpha_1 \sigma_x^2 \mu_y}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \mu \frac{\alpha_2^2 \sigma_y^2 - \alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\alpha_1 = \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\alpha_2 = \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

notice that if $\beta = \frac{1}{2}$ then at $t = 0$ we have a solution iff σ are the same at $t = 0$,

$$\alpha_0 = 0$$

$$\alpha_1 = \frac{1}{2}$$

$$\alpha_2 = \frac{1}{2}$$

$$\alpha_0 = \frac{\alpha_2^2 \sigma_y^2 \mu_x - \alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\alpha_1 = \frac{\alpha_1^2 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

$$\alpha_2 = \frac{\alpha_2 \alpha_1 \sigma_x^2}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

and

$$\mathbb{E} [\hat{p}|p] = \frac{\mathbb{E} \left[\frac{\alpha_1 p + \alpha_2 p_S - \alpha_2 \mu_y}{\alpha_1} \right] \alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \frac{\alpha_1^2 \sigma_x^2 p + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2}$$

and

$$\mathbb{E} [\hat{p}] = \frac{\alpha_1^2 \sigma_x^2 \mu_x + \alpha_2^2 \sigma_y^2 \mu_x}{\alpha_1^2 \sigma_x^2 + \alpha_2^2 \sigma_y^2} = \mu_x$$

Welfare Expressions

The general formula is in the mathematica file, under the restriction $\mu_x = \mu_y$ and $\sigma_x = \sigma_y$ we get

$$W_E = - \left[(\alpha_0 - (1 - \alpha_1 - \alpha_2) \mu)^2 + \left((1 - \alpha_1)^2 + \alpha_2^2 \right) \sigma^2 \right]$$

We have welfare at $t = 0$, where $\mu_S = 1$. Namely

$$W_E = -\beta^2 \left(\underbrace{(\mu_x - \mu_y)^2}_{\text{Prior Bias}} + \sigma_x^2 + \sigma_y^2 \right)$$

$$W_P = -\beta^2 \frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 (1 - \beta)^2 + \sigma_y^2 \beta^2}$$

Then,

$$\frac{W_E}{W_P} = \frac{(\mu_x - \mu_y)^2 + \sigma_x^2 + \sigma_y^2}{\frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 (1 - \beta)^2 + \sigma_y^2 \beta^2}}$$

$$\text{Assume equal } \sigma = \frac{(\mu_x - \mu_y)^2 + 2\sigma^2}{(1 - 2\beta + 2\beta^2)} = \underbrace{\frac{(\mu_x - \mu_y)^2}{\sigma^2}}_{>0} + 2 \left(1 - 2\beta + 2\beta^2 \right)^2 > 2 \left(\frac{1}{2} \right) = 1$$

so if $\mu_x = \mu_y$ (no prior bias), then $W_E(0) = W_P(0)$ iff $\sigma_x = \sigma_y$.

Results

$W_E > W_P$ this is because the bias/variance decomposition

$$\begin{aligned}
 W &= -\mathbb{E} \left[(p - \bar{p})^2 \right] = -\mathbb{E} \left[((1 - \mu_T - \mu_A [2\zeta_A - 1] - \mu_S (1 - \beta)) p + \beta \mu_S p_S + \mu_A (1 - \zeta_A))^2 \right] \\
 &= -\mathbb{E} \left[(p - \hat{p} + \hat{p} - \bar{p})^2 \right] = - \left(\mathbb{E} \left[(p - \hat{p})^2 \right] + \mathbb{E} \left[(\hat{p} - \bar{p})^2 \right] + \cancel{2\mathbb{E} [(p - \hat{p})(\hat{p} - \bar{p})]} \right) \\
 &= - \left(\underbrace{\mathbb{E} \left[(p - \hat{p})^2 \right]}_{\text{Precision of election}} + \underbrace{\mathbb{E} \left[(\hat{p} - \bar{p})^2 \right]}_{\text{Bias of elections}} \right)
 \end{aligned}$$

finally holds, the election have a bias.

The full characterization of the derivative (assuming equal μ and σ)

$$\frac{d}{dt} W_E|_{t=0} = -4\beta\lambda_1\sigma^2 (1 - \beta - \zeta_A)$$

Instead assuming only equal μ we have

$$\frac{d}{dt} W_E|_{t=0} = -2\beta\lambda_1 \left(\beta\sigma_y^2 - \sigma_x^2 (2(1 - \zeta_A) - \beta) \right)$$

so

$$\frac{d}{dt} W_E|_{t=0} > 0 \iff 1 - \beta < \zeta_A$$

or in general

$$\frac{\beta}{2(1 - \zeta_A) - \beta} < \frac{\sigma_x^2}{\sigma_y^2}$$

a sensible condition. Also, λ_1 magnifies either the positive or the negative change local to 0 and in particular if $1 - \beta > \zeta_A$ then more λ_1 is bad for welfare local to $t = 0$. To the contrary,

$$\frac{d}{dt} W_P|_{t=0} = \frac{2(1 - \beta) \beta^2 \lambda_1 \sigma^2 (2\zeta_A - 1)}{\text{sthg}^2} > 0$$

and the welfare of the unconstrained principal is [but this is just a conjecture not falsified by Math plots] always increasing in both λ_1, t .

Conjecture

W_P is increasing in t (and λ_1)— We show that

$$\begin{aligned} \frac{d}{dt} W_P &\propto - \left[\underbrace{2(1-\zeta_A)\mu_S}_{+} \underbrace{\frac{d}{dt}\mu_A}_{?} + \underbrace{(1-2(1-\zeta_A)\mu_A)}_{+} \underbrace{\frac{d}{dt}\mu_S}_{-} \right] \\ &= \underbrace{\exp\{-(\lambda_1+\lambda_2)t\}}_{+} \lambda_1 (2(1-\zeta_A) - \exp\{\lambda_2 t\}) 2(1-\zeta_A) - \exp\{\lambda_2 t\} < 2(1-\zeta_A) - 1 \\ &= 1 - 2\zeta_A < 0 \end{aligned}$$

when computed in

$$\frac{d}{d\lambda_1} W_P \propto - \left[\underbrace{2(1-\zeta_A)\mu_S}_{+} \underbrace{\frac{d}{d\lambda_1}\mu_A}_{?} + \underbrace{(1-2(1-\zeta_A)\mu_A)}_{+} \underbrace{\frac{d}{d\lambda_1}\mu_S}_{-} \right]$$

which has the same sign as

$$\frac{d}{d\lambda_1} W_P \propto - \exp\{\lambda_2 t\} \lambda_2^2 t + 2\lambda_2 [(\exp\{\lambda_2 t\} \lambda_1 t) + (1-\zeta_A) \exp\{(\lambda_2 - \lambda_1)t\} - (1-\zeta_A)(1+\lambda_1 t)] - \dots$$

Furthermore, analyse ζ_A :

$$\begin{aligned} \zeta_A &= - \exp\{\lambda_2 t\} \lambda_2^2 t + 2\lambda_2 [(\exp\{\lambda_2 t\} \lambda_1 t)] - \lambda_1^2 t (\exp\{\lambda_2 t\}) \\ &= -t \exp\{\lambda_2 t\} (\lambda_2 - \lambda_1)^2 + 2\lambda_2 (1-\zeta_A) [\exp\{(\lambda_2 - \lambda_1)t\} - (1+\lambda_1 t)] + 2\lambda_1^2 t (1-\zeta_A) \\ &= \underbrace{-t \exp\{\lambda_2 t\} (\lambda_2 - \lambda_1)^2}_{\text{negative}} + 2\lambda_2 (1-\zeta_A) [\exp\{(\lambda_2 - \lambda_1)t\} - (1+\lambda_1 t) + 2\lambda_1^2 t] \end{aligned}$$

Now if the second addendum is negative then we are done; so assume it is positive, that is

$$\exp\{(\lambda_2 - \lambda_1)t\} - (1+\lambda_1 t) + 2\lambda_1^2 t > 0$$

then the sum is smaller than

$$\begin{aligned}
\underbrace{-t \exp \{\lambda_2 t\} (\lambda_2 - \lambda_1)^2}_{<0} + \lambda_2 \left[\exp \{(\lambda_2 - \lambda_1) t\} (1 + \lambda_1 t) + 2\lambda_1^2 t \right] \\
&= \lambda_1^2 t - \exp \{\lambda_2 t\} (\lambda_2 - \lambda_1)^2 t + \lambda_2 \\
&= \lambda_1^2 t + \lambda_2 \exp \{(\lambda_2 - \lambda_1) t\} - \left[\exp \{
\end{aligned}$$

so it remains to show that this is always negative; if $\lambda_1 \approx 0$

$$-\lambda_2 (1 - \exp \{\lambda_2 t\} (1 - \lambda_2 t)) < -\lambda_2^2 t < 0$$

W_E has interesting comparative statistics due to the interaction with bias. In particular, it seems that for $\beta > \text{stgh}$, then [if there is no prior bias, $\mu_x = \mu_y$] there is a time t such that $Bias(t) = 0$ because the evolution of μ_S, μ_A is such that $\alpha^E = \alpha^P$. This seems interesting, possibly a result to put in a proposition.

Based on our model, we can draw three main results and one additional interesting result.

Inefficiency of twitter economy and non-monotonicity in election times. The first result relates to the political institutions of the digital economy. From our model is that a Twitter-Facebook economy where everyone can speak their mind is not necessarily good: indeed we want only those that went through critical thinking to vote. Following this point, the naturally arising question, *when do we want to hold elections?* Our model clearly implies non-monotonicity in time for election periods.

Typology of voting-users and adverse selection. The second result relates to the typology of voting-users. The “clients” of news outlets, in a micro-foundation of the λ functions are either low i partisans (that look at it for fun) or frustrated critical thinking voting-users that look for some facts (positive predictions). On a related but different point, we can identify the *adverse selection in the vote-force* (under some conditions the strengths of the stereotype pool weakens), and how the format amplifies / reduces this issue (always true that it is better if only types vote, at least in the symmetric case).

Partisan format and compensation effect. The third and most intriguing result relates to the format. We can study the *impact of different storytelling formats* (more in

depth, helps the high i , but how it correlates with α): more in depth, but keep it somehow primitive. In particular and more interestingly, we can allow for *asymmetries*: either there is a “better” policy (say $\beta = 1$, so upon reflecting everyone agrees 1 is right), or stereotypes of one side are less likely to enter critical thinking (evidence that conservatives are overconfident), how does this change the outcome, as well as the incentives for the critical thinking agents (that may vote for those that are less confident because of the bias in the type pool). The problem of asymmetries is that a partisan format, or is the fact that one stereotype is more attractive than the other to make the problem of agents in critical thinking more problematic: remember they are smart but unwise, so they cannot ignore the fact of a stereotyped partisan pool, either because stereotypes are more resistant, or because the shifts the stereotypes. Hence, we propose to explain such a situation by an effect that we label the “Compensation Effect”: When you perceive the device to be partisan in one direction, you vote in the opposite direction when in critical thinking.

The benefits of making voting costly. A resulting and potentially controversial consequence of such an asymmetry is that *voting costs* in this situation may be positive because they can also exclude the strategic types that recognize the stereotype pool is partisan and cannot morally abstain or vote against their type. They can use the excuse not to go to vote.

C.3 Proofs

D Experiment With A Three Cognitive-State Model

In the experiment, we gathered data to decompose the critical thinking process into three stages: S, A, T . Here, S remains unchanged. Now, A denotes an intermediate, transitory stage during which agents experience internal uncertainty regarding the formation of their reasoned preferences. Finally, T denotes the stage in which agents have completed the critical thinking process and have formed their reasoned preferences.

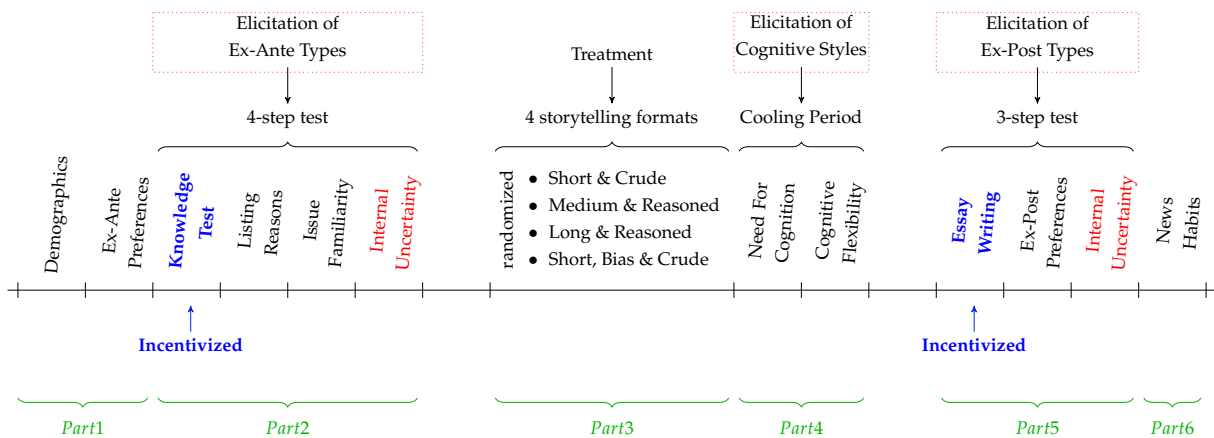


Table 19 shows the classification strategy of participants as Stereotype, Aware, and Type.

Treatment	T	A	S
BEFORE	Knowledge Test Score $> \tau_{KTS}$ Issue Familiarity = 1 Internal Uncertainty $\neq 0$ Reasons List $> \tau_{RL}$	Knowledge Test Score $> \tau_{KTS}$ Issue Familiarity = 1	
AFTER	Psychologists Grade = Pass	Else	

Table 19: CLASSIFICATION STRATEGY BEFORE/AFTER TREATMENT

The analysis presents the frequencies of the three states of participants before and after the treatment

	S_1	A_1	T_1
S_0	475	153	39
A_0	0	21	2
T_0	0	0	30

Table 20: TABLE: FREQUENCIES BEFORE/AFTER TREATMENT