

Social Media, News Consumption, and Polarization: Evidence from a Field Experiment

Ro'ee Levy*

December 19, 2020

Abstract

Does the consumption of ideologically congruent news on social media exacerbate polarization? I estimate the effects of social media news exposure by conducting a large field experiment randomly offering participants subscriptions to conservative or liberal news outlets on Facebook. I collect data on the causal chain of media effects: subscriptions to outlets, exposure to news on Facebook, visits to online news sites, and sharing of posts, as well as changes in political opinions and attitudes. Four main findings emerge. First, random variation in exposure to news on social media substantially affects the slant of news sites individuals visit. Second, exposure to counter-attitudinal news decreases negative attitudes toward the opposing political party. Third, in contrast to the effect on attitudes, I find no evidence that the political leanings of news outlets affect political opinions. Fourth, Facebook's algorithm is less likely to supply individuals with posts from counter-attitudinal outlets, conditional on individuals subscribing to them. Together, the results suggest that social media algorithms may limit exposure to counter-attitudinal news and thus increase polarization.

JEL Codes: D72, L82, L86, O33

*Massachusetts Institute of Technology, roeelevy@mit.edu. I am deeply grateful to Ebonya Washington, Joseph Shapiro, and Dean Karlan for their guidance and support throughout this project. I thank the editor and three anonymous referees for their insightful comments. I also thank Hunt Allcot, Eran Amsalem, Ian Ball, Dirk Bergemann, Leonardo Bursztyn, Alex Coppock, Oren Danieli, Eduardo Fraga, Matthew Gentzkow, Matthew Graham, Eoin McGuirk, Alexey Makarin, Martin Mattsson, Brendan J. Nyhan, Rohini Pande, David Rand, Oren Sarig, David Schönholzer, Katherine Wagner, Jaya Y. Wen, and Mor Zoran, along with seminar participants at Bar Ilan, Berkeley, Booth, Cornell, IDC, Haifa, Hebrew University, Microsoft Research, NBER Digitization, Northwestern, NYU, Tel Aviv, UCLA Anderson, UCSD, Yale, and the ZEW conference on the economics of IT, for their valuable comments. Financial support from the ISPS Field Experiment's Initiative, the Tobin Center for Economic Policy, the Yale Department of Economics, the Yale Program in Applied Economics and Policy, the Yale School of Management, and David Rand is greatly appreciated. All errors are my own. The experiment is registered at the AEA RCT registry (ID 0002713).

In 2019, more than 70% of American adults consumed news on social media, compared to fewer than one in eight Americans in 2008.¹ Based on Pew surveys, Facebook is the dominant social media platform for news consumption, and “among millennials, Facebook is far and away the most common source for news about government and politics” (PEW, 2015, p. 8). As social media becomes a major news source, there are growing concerns that individuals are exposed to more pro-attitudinal news, defined as news matching their ideology, and as a result, polarization increases (Sunstein, 2017).

In this paper, I test whether these concerns are warranted. I analyze the effects of exposure to pro- and counter-attitudinal news outlets by conducting a large online field experiment randomizing exposure to news on Facebook, and by collecting survey, browsing, and social media data.

To motivate the experiment, I first provide descriptive statistics on online news consumption. I show that news sites visited through social media, and specifically Facebook, are associated with more segregated, pro-attitudinal, and extreme news, compared to other news sites visited.

I recruited American Facebook users to the experiment using Facebook ads. After completing a baseline survey, participants were randomly assigned to a liberal treatment, a conservative treatment, or a control group. Participants in the liberal and conservative treatments were asked to subscribe to up to four liberal or conservative outlets on Facebook, respectively (e.g., MSNBC or Fox News), by clicking a “Like Page” button embedded at the end of the survey.² Remarkably, in each treatment, approximately half the participants complied by subscribing to at least one outlet. When individuals subscribe to an outlet on Facebook, posts shared by the outlet may subsequently appear in their Facebook feed. A post usually contains the story’s headline and often includes a link to the full news story on the outlet’s website.

I designed the experiment to have high external validity. A nudge offering subscriptions to outlets is very common on social media and participants could have subscribed to any of these outlets, at no cost, without the intervention. Besides the offer, the experiment did not directly intervene in any behavior. The news supplied to participants was the actual news provided by leading media outlets during the study period. Facebook’s algorithm determined which of the posts shared by the outlets appeared in the participants’ Facebook feeds. Finally, participants decided whether to skip, read, click, or share posts. As a result, the effect of the intervention is almost identical to the experience of millions of Americans who subscribe to popular news outlets on Facebook.

I estimate the effect of the intervention on exposure to news in the Facebook feed, news sites visited, news shared, political opinions, and affective polarization, defined as negative attitudes toward the opposing political party. Affective polarization is a primary outcome of interest since this measure of polarization has been increasing (Iyengar and Krupenkin, 2018), and there are concerns over its implications for governance, accountability of elected officials, and even labor markets (Iyengar et al., 2019).

To measure subscriptions to outlets on Facebook and posts shared, I asked participants to log in to the survey using their Facebook account. To measure exposure to news in the Facebook feed and visits to news sites, I developed a Google Chrome extension and asked a subset of participants who took the survey on a computer using Chrome to install it. To estimate the effect on opinions and attitudes, I invited participants to an endline survey approximately two months after the intervention. My sample is composed of 37,494 participants who completed the baseline survey. 34,592 of those participants provided access to the posts they shared for at least two weeks, 1,835 installed the extension for at least two weeks, and 17,635 took the endline survey.

This paper has four main findings. First, exposure to news on social media substantially affects online news consumption. Following increased exposure to posts from the randomly offered outlets, participants visited the websites of the outlets, even when the outlets did not match their ideology. Visiting the websites had a substantial effect on the mean slant of participants’ overall online news consumption. The difference between the intention-to-treat (ITT) effects of the liberal and conservative treatments on the slant of news

¹2008 figure is based on the Pew Research Center 2008 Biennial Media Consumption Survey. The 2019 figure is based on the Pew Research Center American Trends Panel Wave 51, July 2019.

²To simplify terminology, throughout the paper I will describe the action of “liking” a page of a news organization as subscribing to an outlet on Facebook.

sites visited in the two weeks following the intervention is 14% of the difference in the slant of sites visited by liberals and conservatives in the control group.

Various economic theories explain why individuals optimally prefer news that matches their ideology (Gentzkow, Shapiro and Stone, 2015). However, I find that news consumption strongly responds to an exogenous shock to the feed, meaning that individuals often consume news incidentally, and do not fully re-optimize their browsing behavior to keep the slant of the news sites they visit constant. The results imply that social media algorithms can substantially alter news consumption habits and that while social media is associated with pro-attitudinal news, individuals are willing to engage with counter-attitudinal news when it is made more accessible on social media.

My second finding is that exposure to counter-attitudinal news *decreases* affective polarization, compared to pro-attitudinal news. I construct an affective polarization index measuring attitudes toward political parties. The index includes questions such as how participants feel toward their own party and the opposing party, i.e., a “feeling thermometer”. When estimating the effects on polarization, I redefine the treatments as pro- and counter-attitudinal. For example, a counter-attitudinal treatment is a liberal treatment assigned to a conservative participant or a conservative treatment assigned to a liberal participant. The ITT and treatment-on-treated (TOT) effects of the counter-attitudinal treatment on the affective polarization index, compared to the pro-attitudinal treatment, are -0.03 and -0.06 standard deviations, respectively. The TOT effect should be interpreted as the effect on individuals who subscribe to new outlets when nudged to subscribe. Comparing each treatment to the control group suggests that the effect on polarization is driven by the counter-attitudinal treatment but this result should be interpreted cautiously since participants in the control group were more likely to complete the endline survey (there is no differential attrition between the two treatment arms).

I compare the results to existing benchmarks by focusing on the feeling thermometer questions. The experiment’s ITT and TOT effects decreased the difference between participants’ feelings toward their own party and the opposing party by 0.58 and 0.96 degrees on a 0-100 scale over two months, respectively. For comparison, based on the American National Election Survey (ANES), this measure of affective polarization increased by 3.83-10.52 degrees between 1996 and 2016.³

Third, in contrast to the effect on attitudes, I do not find evidence that the slant of news outlets affects political opinions. The effect of the liberal and conservative treatments on a political opinions index focusing on issues and political figures covered during the study period is small in magnitude, precisely estimated, and not statistically significant.

The paper’s fourth finding is that Facebook’s algorithm may limit exposure to counter-attitudinal news. I show that participants in the counter-attitudinal treatment were exposed to substantially fewer posts from the outlets they subscribed to in the intervention, compared to participants in the pro-attitudinal treatment.

Combined, the results paint a complicated picture. On the one hand, Facebook’s algorithm seems to filter counter-attitudinal news, probably since it attempts to personalize news based on the user’s behavior and perceived interests. While it is not possible to estimate the effect of specific posts filtered by the algorithm, I show that exposure to counter-attitudinal news decreases affective polarization. This suggests that social media algorithms may be increasing polarization. On the other hand, this paper also shows that individuals are willing to engage with counter-attitudinal news, and social media platforms provide a setting where a subtle nudge can substantially diversify news consumption and consequently decrease affective polarization.

This paper contributes to the literature on social media and news consumption. In his seminal book “The Filter Bubble,” Eli Pariser warned that the “era of personalization is here” (Pariser, 2011, p. 19). However, recent reviews concluded that “we lack convincing evidence of algorithmic filter bubbles in politics” (Guess et al., 2018, p. 12). Papers in this literature typically estimate segregation in online news based on cross-sectional analysis of browsing behavior (Gentzkow and Shapiro, 2011; Flaxman, Sharad and Rao, 2016;

³The increase in polarization depends on the weights and the respondents included in the sample. When using the ANES face-to-face sample for consistency (Boxell, Gentzkow and Shapiro, 2018), the increase is 3.83. When including also the 2016 web sample (Iyengar et al., 2019), the increase is 10.52. The ANES top-codes the thermometer at 97 degrees. The results are almost exactly the same when I top-code the results.

Peterson, Shared and Iyengar, 2019; Guess, 2020). Since they lack social media data, these papers cannot measure segregation *within* one’s social media feed. One exception is a paper analyzing Facebook data, arguing that exposure to counter-attitudinal news shared by *friends* is mostly limited by individual choices and not by algorithmic ranking (Bakshy, Messing and Adamic, 2015). The paper analyzes large data but does not exploit exogenous variation. I advance the literature by generating experimental variation in subscriptions to outlets and collecting data on exposure to posts from those outlets. This allows me to decompose the mechanisms limiting exposure to counter-attitudinal news and demonstrate the existence of a filter bubble, i.e., that Facebook’s algorithm is more likely to expose individuals to news matching their ideology, conditional on subscription.

My findings contribute to the literature on social media and polarization by generating variation in the main mechanism through which social media is suspected to increase polarization: the distance between individuals’ ideology and the slant of their news consumption. Related papers show that the Internet and Facebook may increase polarization (Lelkes, Sood and Iyengar, 2015; Allcott et al., 2020), but based on demographics, they may not be the primary driver in the rise of polarization (Boxell, Gentzkow and Shapiro, 2018).⁴ These papers focus on the reduced-form effect of social media and do not identify the causal effect of pro- or counter-attitudinal news. Indeed, a recent review argues that “it is far from clear ... that partisan news actually causes affective polarization” (Iyengar et al., 2019, p. 135). To the best of my knowledge, this paper provides the first experimental evidence that counter-attitudinal news decreases affective polarization and thus demonstrates that nudges diversifying social media news exposure can be effective.

This study also contributes to a well-established literature on media persuasion by randomly assigning subscriptions to news outlets. Survey experiments (e.g., Coppock, Ekins and Kirby 2018) and papers with quasi-experimental designs (e.g., DellaVigna and Kaplan 2007) find that individuals are persuaded by the news they consume.⁵ While in many contexts field experiments are considered the gold standard for estimating causal effects, field experiments estimating media effects are not common. One notable exception is a study randomizing subscriptions to the Washington Post and Washington Times, which does not find an effect on opinions but is limited by a relatively small sample size (Gerber, Karlan and Bergan, 2009). This paper studies a different setting, social media, and shows how the unique features of this setting affect news exposure. Focusing on social media also allows me to analyze engagement with news and quantify the effect of news exposure.

Methodologically, this paper contributes to a growing literature conducting online media-related experiments (Bail et al., 2018; Jo, 2018; Chen and Yang, 2019; Mosquera, Odunowo and Mcnamara, 2019; Allcott et al., 2020) by demonstrating how an experiment can exploit social media’s existing infrastructure to gradually distribute news to participants in a natural setting. In contrast to most online experiments, participants were not asked to consume any content or continue complying with the treatment over time, nor did they receive frequent reminders of the experiment. The natural, unobtrusive intervention means that it is unlikely that experimenter effects drive the study’s result. To precisely detect the small effects that are expected as a result of a subtle intervention, I collect a sample size that is an order of magnitude larger than most other related experiments.

I Background: Facebook

This study focuses on Facebook since it is the dominant social media platform, used by seven out of ten American adults. Most of these users visit Facebook several times a day,⁶ and the platform accounts for 45% of all time spent on social media (Williamson, 2018). Despite its prominence, Facebook has been understudied, especially compared to Twitter (Guess et al., 2018).

⁴Other studies estimating the effect of social media on political behavior include Bursztyn et al. (2019), Müller and Schwarz (2019), and Enikolopov, Makarin and Petrova (2020). See Zhuravskaya, Petrova and Enikolopov (2020) for a recent review.

⁵Other studies estimating media effects on political outcomes include Chiang and Knight (2011), Gentzkow, Shapiro and Sinkinson (2011), Durante, Pinotti and Tesei (2019), and Okuyama (2019). See Strömberg (2015) for a review.

⁶Facebook usage is based on the Pew Research Center January 2019 Core Trends Survey.

The most important Facebook feature is the news feed, where users scroll through a list of posts curated by Facebook’s algorithm. Posts in the feed are typically shared by the user’s Facebook friends, shared by Facebook pages the user subscribes to (“likes”), or are sponsored (advertisements shared by pages to promote content). The posts may include text, video, pictures, and links.

Facebook is a very popular source for news consumption. In 2019, 52% of Americans reported getting at least some of their news on Facebook, more than the share of Americans getting news on all other social media platforms combined.⁷ While this study focuses on the US, understanding Facebook’s influence has global implications. A Reuters Institute survey found that in 37 out of 38 middle and high-income countries surveyed, more than 20% of the population consumed news through Facebook weekly (Reuters Institute, 2019). A paper analyzing the survey’s data concluded that Facebook “reaches the widest international audience of any media organization in our sample” (Kennedy and Prat, 2019, p. 10).

With Facebook’s growing influence, it has faced several controversies in recent years, including an effort by the Russian-based Internet Research Agency to influence the elections, the spread of fake news during the 2016 US election cycle, and Cambridge Analytica’s attempt to assist campaigns with personally targeted ads. The concerns over each of these scandals were based on the assumption that individuals are easily persuaded by political content on social media.

II Design and Data

This section summarizes the experimental design, data, and empirical strategy. The design of the experiment is also presented in Figure 1.

A Experimental Design

I recruited American adults to the experiment in February-March 2018 using Facebook ads. 978,628 people saw the ads, 87,648 people clicked the links in the ads, and approximately half of those began the survey. For more details on the ads see Appendix A.1.1. Individuals who clicked the ads were directed to the survey landing page, where they reviewed the consent form and began the survey by logging in using their Facebook account.

After logging in, and before treatment assignment, four *potential* liberal outlets and four *potential* conservative outlets were defined for each participant. The same eight potential outlets were defined for each participant unless a participant already subscribed to one of the outlets in baseline, in which case it was replaced with an alternative outlet, to ensure only new outlets would be offered. Toward the end of the survey, participants were randomly assigned to a liberal treatment, a conservative treatment, or a control group, with the randomization blocked by participants’ self-reported baseline ideology.⁸ Participants in the conservative treatment were offered to subscribe (“like”) to their four potential conservative outlets and participants in the liberal treatment were offered to subscribe to their four potential liberal outlets. Participants in the control group were not offered any outlets.

I nudged participants to subscribe to the outlets by explaining that subscribing could expose them to new perspectives. Participants were not required to subscribe to any outlet and did not receive monetary compensation for subscribing. The intervention did not provide exclusive access to these outlets, and any individual can subscribe to these outlets on Facebook, regardless of the intervention. Since participants were logged into their Facebook account when taking the survey, the offer to subscribe was integrated within the survey, and the only action required by participants was to click the standard Like Page button. Facebook

⁷Calculation based on the Pew Research Center American Trends Panel Wave 51.

⁸Respondents were asked where they position themselves on a 7-point ideological scale, with an additional option of “I haven’t thought about it much.” Each block is composed of three sequential participants who chose the same answer. The first participant in a block was randomly assigned to one of the three treatment groups, the second participant was randomly assigned to one of the two remaining groups, and the third participant was assigned to the remaining group.

users often encounter this button, for example when Facebook suggests pages they may be interested in or when outlets promote their page. Appendix Figure A.1 provides an example of the intervention.

After participants subscribed to an outlet by “liking” its Facebook page, posts from the outlet appeared in their feeds, among many other posts, according to Facebook’s algorithm. Participants decided whether to read a post, click a link, share a post or unsubscribe from an outlet, just like the decisions they make regarding other posts appearing in their feed. Due to the simple common intervention, the organic nature of any subsequent effect, and the fact that participants were not reminded of the intervention, experimenter effects are unlikely to play a large role in explaining the effects, at least compared to similar studies.⁹ Because individuals can subscribe to outlets on Facebook at no cost and no monetary incentives were provided, the intervention is scalable.

B The Setting: Media Outlets and the News Environment

The primary liberal outlets offered in the experiment were HuffPost, MSNBC, The New York Times, and Slate. The primary conservative offered outlets were Fox News, The National Review, The Wall Street Journal, and The Washington Times. The news outlets were chosen to ensure participants are offered a diverse set of popular outlets (Fox News and the New York Times are two of the three most popular news pages on Facebook) with a clear ideological slant. Appendix Table A.1 displays the full list of primary outlets offered, along with the alternative liberal and conservative outlets.

Figure 2 shows that the men and women mentioned most often in posts shared by the eight primary outlets and the two main alternative outlets are political figures. Unsurprisingly, President Trump is the dominant figure mentioned. Political stories that made headlines during the study period can be observed in the figure: Trump’s alleged affair with Stormy Daniels, Robert Mueller’s investigation, and the negotiation with North Korea’s leader, Kim Jong Un. The figure also demonstrates that liberal outlets focused on scandals related to the presidency and mentioned Michael Cohen, Stormy Daniels, Scott Pruitt, and Vladimir Putin, more often than conservative outlets.

C Data Collection and Subsamples

C.1 Experiment Data

The analysis of the experiment relies on three datasets: self-reported survey data, Facebook data, and browser data. This is among the first studies combining experimental variation with social media and news-related browsing data. Table 1 presents the main datasets and subsamples analyzed.

Survey Data

The endline survey measures self-reported political opinions, affective polarization, and changes in news consumption habits. 17,635 participants took the endline survey and constitute the *endline survey subsample*.

Facebook Data on Pages Liked and Posts Shared

Participants logged in to the survey using their Facebook account, through a Facebook app created for the project. They were asked to provide separate permissions to access the pages they subscribe to and posts they share. Providing permissions was voluntary, they could be revoked at any time, and were revoked automatically approximately two months after participants logged in to a survey. I observe all posts shared or pages liked until permissions are revoked. Since baseline subscriptions were required to define the potential outlets, participants who did not provide initial permissions to access their subscriptions are excluded from the baseline sample.¹⁰

⁹Participants were asked at the end of the survey what they think is the purpose of the study. Appendix C.1 shows that participants understood the study was about media and politics and that there do not appear to be dramatic differences between the answers of participants in the pro- and counter-attitudinal treatments.

¹⁰While providing permissions was not required to complete the survey or to be eligible for any rewards, the vast majority of participants who completed the survey provided these permissions. Participants who revoked permissions after the intervention are included in the baseline sample.

Data on posts shared is used to estimate the effect of the intervention on political behavior. I exclude posts sharing photos, albums, music, and events. The remaining posts typically include text with a link or an embedded video. Since posts shared are observable to the participant's social network or the general public, sharing posts can have a direct cost to the reputation of the participant. Approximately 92% of baseline participants provided access to the posts they shared for at least two full weeks following the intervention constituting the *access posts subsample*.

Extension Data on Browser Behavior and the Facebook Feed

Participants who completed the baseline survey using Google Chrome on a computer were asked to install a browser extension collecting data on the Facebook feed and news-related browsing behavior, in exchange for a small reward. The offer was made toward the end of the survey, but before the intervention, to ensure take-up is not affected by the intervention. 2,262 of the 8,084 participants who were offered the extension, installed it. I focus on 1,835 participants who kept the extension installed for at least two weeks and constitute the *extension subsample*.

The Facebook feed data is used to analyze news exposure by estimating how often participants were exposed to posts from outlets on Facebook. I observe the posts that participants saw when they used their computer mouse to scroll their feed. I do not observe whether a post is a sponsored advertisement, but identify suspected ads as posts in the feed from pages participants did not subscribe to and posts appearing in the feed repeatedly. I attribute a post to a news outlet if it was created by the outlet's Facebook page or contains a link to the outlet's domain.¹¹ While the variation generated by the experiment is in subscriptions to the outlets' Facebook pages, my analysis includes news articles shared by the participants' friends, to accurately capture total exposure to news outlets on Facebook.

The browsing behavior data is used to estimate the effect on the news sites participants visited. The extension can greatly reduce measurement error, compared to self-reported estimates of news consumption, as individuals' self-reported media habits may be more polarized than their actual news consumption (Guess, Nyhan and Reifler, 2017).

The extension data was only collected when participants used a computer while being signed into their Chrome account. In practice, individuals often use Facebook and browse news sites on a mobile device or at work, where they may use a different browser. Therefore, the estimates for the number of posts participants were exposed to in their feed and the number of sites they visited are lower bounds.¹²

Additional details on the survey, Facebook, and extension data can be found in Appendices A.1, A.2, and A.3, respectively.

Subsamples

The datasets define three separate subsamples. To maximize power, when analyzing the effects on opinions and attitudes, I focus on the *endline survey subsample* and when analyzing media outcomes, I focus on the *extension subsample* and the *access posts subsample* (or their overlap). Appendix Table A.2 presents descriptive statistics on the subsamples and shows that the extension subsample is more liberal and older, as would be expected when excluding participants who took the survey on a smartphone. The share of compliers is greater in the extension subsample, which assists in detecting treatment effects despite the smaller sample size.

C.2 External Data

Outlets

¹¹To match URLs with news outlets, I first convert over ten million URLs to their final endpoint, following redirects. This is required since many links on Facebook are based on URL-shortening services such as tinyurl.com.

¹²In the baseline survey, participants were asked how many links to articles about government and politics they clicked on Facebook in the past 24 hours using a computer and on a mobile phone. Among participants in the extension subsample who provided a numerical answer under 1,000, approximately 72% of news links were clicked on a computer, so it is likely that most, but not all data is collected for these participants.

I measure the slant of news at the outlet level, the common method used in the literature. I determine an outlet's slant according to a dataset by Bakshy, Messing and Adamic (2015) defining the slant of 500 news domains based on the self-reported ideology of Facebook users sharing articles from the domains. Using this definition, a completely liberal outlet has a slant of approximately negative one, a middle-of-the-road outlet has a slant of approximately zero, and a completely conservative outlet has a slant of approximately one. I use this measure of slant since it is based on 2014 data, and thus more recent than other common measures, and since it covers a large number of online news outlets. The dataset correlates well with other measures of slant (e.g., Gentzkow and Shapiro, 2010). I refer to outlets in this dataset as *leading news outlets*. I exclude from the dataset several popular domains which are clearly not news outlets or that serve mostly as portals, and merge several outlets that are associated with multiple domains. I manually determine the Facebook pages of leading outlets by searching for pages with names similar to each outlet's domain. Facebook pages were found for 370 outlets.

Comscore Browsing Data

To provide descriptive statistics on news consumption outside the experimental sample, I analyze the 2017 and 2018 Comscore WRDS Web Behavior Database Panel (Comscore, 2018). Each observation in the dataset is a domain visited by a specific computer along with the referral domain. I merge this dataset with the list of leading news outlets. The Comscore data provides several advantages. The combined 2017 and 2018 samples include 94,342 individuals who visited at least one news site. Previous studies have shown that the panel is representative of online buyers in the United States (Hortacsu, Wildenbeest and De Los Santos, 2012). Finally, the data has been collected for previous years and used by other researchers (Gentzkow and Shapiro, 2011), allowing me to estimate changes in news consumption over time. I classify the channels through which visitors reached websites as social, search, or direct visits. Facebook is by far the dominant referral source in the social category.

For more details on the outlet and Comscore datasets, see Appendices A.4 and A.5, respectively.

D Outcomes

D.1 Media

I measure subscriptions to outlets on Facebook, exposure to news in the Facebook feed, news sites visited, and posts shared, using the following outcome measures. First, I estimate the direct effect of the experiment according to the number of times participants engaged with the *potential outlets* (the four liberal outlets and the four conservative outlets defined for each participant). For example, I measure the number of posts participants observed from their potential liberal and conservative outlets in their feed. Second, I measure the mean slant of all *leading news outlets* participants engaged with. Third, to measure the effects of the pro- and counter-attitudinal treatments on total news consumption, I define a *congruence scale*, calculated as the mean slant of news consumed, multiplied by (-1) for liberal participants. This scale has a higher value when individuals consume more extreme news matching their ideology. Fourth, I measure the *share of counter-attitudinal news*, defined as the share of news from counter-attitudinal outlets among all news from pro- and counter-attitudinal outlets.

D.2 Opinions and Attitudes

I analyze the effects of news exposure on two primary outcomes: political opinions and affective polarization. For both outcomes, an index is composed by taking an average of all the valid non-missing index components and then standardized by subtracting the control group mean and dividing by the control group's standard deviation.

The political opinions index is composed of twenty survey questions focusing on domestic political issues and political figures covered in the news during the study period, such as new tariffs, the March For Our Lives Movement, and the investigation regarding Russian interference in the elections. Each outcome

variable is defined such that a higher value is associated with a more conservative opinion and then standardized.

The affective polarization index is composed of five outcomes. First, I use the feeling thermometer questions (*feeling thermometer*). Second, participants were asked how well the following statement describes them on a scale from 1 to 5: “I find it difficult to see things from Democrats’/Republicans’ point of view” (*difficult perspective*). Third, participants were asked a similar question on the following statement: “I think it is important to consider the perspective of Democrats/Republicans” (*consider perspective*).¹³ Fourth, participants were asked if they think the Democratic and Republican parties have a lot (3), some (2), a few (1), or almost no good ideas (0) (*party ideas*). For each of the four previous measures, I calculate the difference between attitudes toward the party associated with the participant’s ideological leaning and attitudes toward the opposing party, a typical measure of affective polarization. Fifth, participants are asked if they would feel very upset (2), somewhat upset (1), or not upset at all (0) if they had a son or daughter who married someone from the opposing party (*marry opposing party*).¹⁴ Each outcome variable is defined such that a higher value is associated with more polarization and then standardized.

E Empirical Strategy

When estimating the effects of the intervention on engagement with the liberal and conservative outlets, the slant of news participants engaged with, and their political opinions, I compare the liberal and conservative treatments. When measuring the effects on polarization and engagement with pro- and counter-attitudinal outlets, it no longer makes sense to use these treatments (a conservative treatment is not expected to make participants more or less polarized than a liberal treatment), and therefore I focus on the pro- and counter-attitudinal treatments. This strategy follows the study’s pre-analysis plan, discussed in Appendix B.2.

Liberal and Conservative Treatments I estimate the following ITT regression:

$$Y_i = \beta_1 T_i^L + \beta_2 T_i^C + \alpha X_i + \varepsilon_i \quad (1)$$

where $T_i^L, T_i^C \in \{0, 1\}$ is whether participant i is assigned to the liberal or conservative treatment, respectively. When estimating the effect on political opinions, I focus on the difference between the liberal and conservative treatments, by testing whether $\beta_1 < \beta_2$ (i.e., the conservative treatment made participants more conservative, compared to the effect of the liberal treatment). To increase power, I control for the following set of covariates, X : self-reported ideology, party affiliation, approval of President Trump, ideological leaning, age, age squared, gender, and baseline questions measuring political opinions that are similar to questions used in the endline survey. Appendix B.3 describes the control variables. When estimating the effect on media outcomes, I only control for the outcomes in baseline, when they exist. All regressions use robust standard errors unless noted otherwise.

Pro-Attitudinal and Counter-Attitudinal Treatments I estimate the following ITT regression:

$$Y_i = \beta_1 T_i^A + \beta_2 T_i^P + \alpha X_i + \varepsilon_i \quad (2)$$

¹³Both statements are based on an empathy index developed by Robb Willer, Jamil Zaki, and Emily Reit, loosely based on the Interpersonal Reactivity Index (Davis, 1980).

¹⁴Participants stating in the endline survey that they are Republicans or Democrats were asked how they would feel if they had a son or daughter who married a Democrat or Republican, respectively. Participants who did not identify with either party were asked about one of the parties randomly. I asked participants only about the opposing party since I was concerned they would find it odd to state how upset they would be if they had a son or daughter who married someone from their own party. However, conditioning the question on an endline variable could potentially bias the result. For example, if some participants were affected by the counter-attitudinal treatment, and as a result, no longer identified with their party, they were less likely to be asked about the opposing party in endline and the average participant asked about the opposing party would be slightly less moderate in this treatment arm. I include this measure in the affective polarization index since it is the only social-distance measure in the index, it is included in the pre-analysis plan, and any bias is expected to go against the direction of my findings. Appendix Table A.13 shows that the results are robust to excluding this measure from the index.

where $T^A \in \{0, 1\}$ is whether the participant was assigned to the counter-attitudinal treatment, defined as a liberal treatment assigned to a conservative participant or a conservative treatment assigned to a liberal participant. $T^P \in \{0, 1\}$ is whether the participant was assigned to the pro-attitudinal treatment, defined as a liberal treatment assigned to a liberal participant or a conservative treatment assigned to a conservative participant. X is the same set of control variables used when analyzing the effect on political opinions, with baseline measures of political opinions replaced with baseline measures of affective polarization. $\beta_1 < \beta_2$ tests whether individuals become more polarized when assigned to pro-attitudinal news, compared to counter-attitudinal news.

I determine whether participants are liberal or conservative (their ideological leaning) according to the party they identify with or lean toward. If participants do not identify with either the Democratic or Republican Party, their ideological leaning is defined according to their self-reported ideology, and if they do not identify as liberal or conservative, it is defined according to the candidate they supported in the 2016 elections.¹⁵

F Balance and Attrition

Table 2 presents descriptive statistics for participants in the baseline sample and shows that the sample is balanced. Appendix Table A.3 presents a balance table according to whether the treatment matched the participant's ideology (pro- or counter-attitudinal), and shows that the sample is balanced along the redefined treatment arms as well. The sample size in this table is slightly smaller because it excludes participants for whom an ideological leaning cannot be defined.

Similar to other opt-in panels, the sample is not nationally representative. Participants tend to be more liberal than the US population and, as expected, more participants say that they get most of their news on social media (18%), compared to the national population (13%). The share of female participants and the average age is similar to the US population. Self-reported exposure to news on Facebook in line with one's views is similar to US Facebook users. Overall, the sample seems at least as representative as samples of Mechanical Turk users (Berinsky, Huber and Lenz, 2012).¹⁶

Table 2 and Appendix Table A.3 also test for differential attrition among the three subsamples. The access posts and extension subsamples have low attrition rates compared to baseline takeup (as shown in Table 1) and very small differences in attrition by treatment arm. Therefore, their results are unlikely to be affected by attrition.¹⁷ However, more participants completed the endline survey in the control group (49%), compared to the liberal (46%) and conservative (46%) treatment arms. The differential attrition mostly stems from a small share of participants in the conservative and liberal treatments who did not complete the final screen of the baseline survey after they encountered the intervention.

Appendix Tables A.4 and A.5 present balance tables for the endline survey subsample and show that despite the attrition, the two treatment arms and control group are similar on observables. Participants in the pro-attitudinal treatment who completed the endline survey are *not* substantially more polarized in baseline than participants in the counter-attitudinal treatment. Moreover, there is no differential attrition between the conservative and liberal treatments and no differential attrition between the pro- and counter-attitudinal treatments. When estimating the effect on the primary endline survey outcomes, I compare the two treatment arms to each other to mitigate concerns over differential attrition. Still, it is possible that attrition could affect the results.

¹⁵Approximately 3% of participants do not identify with the Republican or Democratic Party, do not self-identify as liberals or conservatives, and did not support Trump or Clinton. They are excluded from the analysis when analyzing the effect of the pro- and counter-attitudinal treatments. The effect on affective polarization is robust to including only participants who identify with or lean toward the Democratic or Republican Party.

¹⁶One advantage of the sample is that Facebook users are not experienced, semi-professional survey takers. Participants were asked in the endline survey how many additional surveys they completed in the past month, the median answer is 1 and the mean answer is 7. For comparison, a 2014 study found that the median Mechanical Turk worker reported participating in 20 academic studies in the week before the question was asked (Rand et al., 2014).

¹⁷There is a very small, but statistically significant difference between the conservative treatment and the other groups in the number of participants who provided permissions to access their posts for two weeks following the intervention (the *Access Post, Two Weeks* variable in Table 2). However, this minimal difference seems to be random, since it already existed before the intervention, as can be seen by the *Access Post, Pre-Treat* variable.

G Compliance

Throughout the analysis, I focus on ITT estimates. To measure the effect of complying with the treatment, defined as subscribing to at least one offered outlet, I also analyze TOT estimators by regressing the dependent variable on compliance with each treatment and instrumenting compliance with the random treatment assignment. Since the intervention only offers new outlets to participants, defiers do not exist in this experiment.¹⁸ Because compliance is defined as liking an outlet when it was offered, always-takers do not exist either.¹⁹ If compliers are more likely to engage with the outlets and be affected by them, perhaps because they are more interested in the content, the TOT is expected to be larger than the average treatment effect.

In the entire baseline sample, 59% of participants who were offered pro-attitudinal outlets complied with the pro-attitudinal treatment and subscribed to at least one outlet, compared to 48% of participants offered counter-attitudinal outlets. Table 3 shows that participants were more likely to subscribe to outlets they are familiar with, to outlets with a perceived ideology similar to their own ideology, and to outlets they perceive as more moderate. Appendix Table A.6 presents descriptive statistics on the compliers by treatment arm and shows that liberals, women, and participants who subscribe to more outlets on Facebook were more likely to comply with both treatments. To test whether participants open to new ideas comply more often with the treatments, I use two questions from a brief measure of the big five personality traits (Gosling, Rentfrow and Swann, 2003), along with self-reported certainty in political opinions, and exposure to counter-attitudinal news in baseline. Based on these measures, compliers with the counter-attitudinal treatment are slightly more open than non-compliers, but the differences are not large (0.12-0.19 standard deviations) and some of these differences exist to a lesser degree when comparing compliers and non-compliers among participants assigned to the pro-attitudinal treatment.

This section deals with immediate compliance with the intervention, which is especially useful when interpreting the TOT effects. However, the experiment is designed to allow participants to opt-out of news content at any stage. They could always unsubscribe from the offered outlets or ignore posts from the outlets appearing in their feed. Therefore, the effects found will probably be driven by participants who decide to consume the content offered when it becomes accessible. This feature increases the external validity of the results because these participants are often the policy-relevant population, as they are more likely to engage with the offered outlets in other circumstances as well.

III Descriptive Analysis: Segregation in Online News Consumption

Is the rise of social media associated with a change in news consumption patterns? In this section, I present descriptive statistics on segregation in social media and online news. I calculate two main measures: isolation and segregation.

Isolation measures whether conservatives and liberals visit different websites. It is defined as the difference between exposure to conservatives in websites visited by conservatives and exposure to conservatives in websites visited by liberals. Exposure to conservatives is the share of conservatives visiting each set of websites. Intuitively, if conservatives tend to visit websites visited by many other conservatives, while liberals tend to visit websites visited by few conservatives, the isolation measure is higher. To make the measure comparable to estimates by Gentzkow and Shapiro (2011), I aggregate visits at the daily level and use the adjusted leave-out estimator of isolation.

¹⁸Defying the experiment would mean unsubscribing from an offered outlet, but participants are only offered outlets they are not already subscribed to. There are rare cases where I only observe a partial list of outlets in baseline and as a result, a participant could have been offered an outlet she already subscribed to and “unliked” the outlet’s page instead of “liking” it. However, I estimate that I observed a partial list of outlets for less than 1% of participants and I do not have evidence that participants unsubscribed from outlets as a result of the intervention.

¹⁹In a handful of cases, participants subscribed to potential outlets even when the outlets were not offered, possibly since the survey included questions about these outlets. However, these cases are extremely rare and therefore, I am not defining them as compliance for simplicity. When focusing on the two weeks following the intervention instead of immediate compliance, an always-taker would be defined as a participant who would subscribe to a potential outlet in that period, regardless of the intervention. In the control group, only 0.2% and 0.5% of participants subscribed to a potential conservative or liberal outlet, respectively, in the two weeks following the intervention.

Segregation is defined as the scaled standard deviation of the slant of news sites visited by participants. To keep this measure in the unit interval, the slant of outlets is normalized to range from zero to one. A higher value means that there is a greater difference in the slant of news consumed by two random individuals. The measures are formally defined in Appendix B.1.

A Segregation in Online News

I find that news consumed through social media is more segregated and extreme than news consumed through other channels. Figure 3 shows that in the 2017-2018 Comscore sample, the segregation index is 0.18 for news sites visited through search engines, 0.21 for news sites visited directly, and 0.28 for news sites visited through social media. I cannot precisely calculate the isolation measure for the Comscore panel since individual ideology is not observed. Instead, in Panel 2 of Figure 4, I analyze isolation based on the extension subsample and show that isolation is greater in news sites visited through Facebook, compared to news sites visited through other means. The analysis is based on participants assigned to the control group and includes data from the first eight weeks after the extension was installed. The full results are presented in Appendix Table A.7, which also shows that total segregation is similar in the extension and Comscore datasets.

The increased segregation for news sites visited through social media could stem from the individuals using social media to consume news. Appendix Table A.8a presents the segregation among the 8,882 individuals in the Comscore sample who visited multiple news sites through Facebook and through other sources. As all the individuals in this group consume news through both channels, the comparison better isolates the effect of the medium. While the share of news sites visited through Facebook is much greater among these individuals (26%), sites visited through Facebook remain substantially more segregated than sites visited through other means.

Figure 5a presents the distribution of the mean slant of news consumption for these individuals and shows that news sites visited through Facebook are more extreme. When visiting news sites through Facebook, 57% of individuals consume news that is on average more conservative than the Wall Street Journal or more liberal than the Washington Post, and when visiting news sites through other sources, 39% of individuals consume such partisan news.²⁰

Figure 5b shows a clear correlation between the consumers' ideology and the slant of their news consumption. I proxy for ideology using the share of donations to Republican candidates in the consumers' zip codes in the 2016 and 2018 election cycles, based on FEC data. The slope for news consumed through Facebook is steeper than the slope for news consumed through other sources, indicating that sites visited through Facebook tend to better match the consumers' ideology.

Has segregation in online news consumption increased? In the extension sample, the *segregation* index for all online news is 0.20 when I define the slant of outlets based on Bakshy, Messing and Adamic (2015) and 0.23 when I define the slant based on the ideological leaning of participants (Peterson, Shored and Iyengar, 2019). These segregation levels are similar to a value of 0.25 found by Peterson, Shored and Iyengar (2019) using 2016 data from the Wakoopa toolbar and substantially larger than a value of 0.11 found by Flaxman, Sharad and Rao (2016) using 2013 Bing toolbar data. To compare the *isolation* index to previous estimates, I use visit-level measures of isolation (row 6 in Appendix Table A.9b), which give more weight to individuals who visit more news sites. The isolation of browsing behavior in the extension sample is 0.22, similar to a value of 0.21-0.24 calculated by Peterson, Shored and Iyengar (2019) and larger than a value of 0.07-0.08 calculated by Gentzkow and Shapiro (2011). One limitation with this comparison is that while I attempt to make the samples comparable, each study still analyzes the data slightly differently. In Appendix Table A.8b, I provide a cleaner comparison of segregation levels by comparing 2007-2008 and 2017-2018 Comscore data and do not find substantial changes in segregation.

²⁰Washington Post and the Wall Street Journal are in the 36th and 63rd percentile of the Bakshy, Messing and Adamic (2015) dataset. When using the 25th and 75th percentile, which are similar to Boston Globe and Fox News, 19% of individuals consume news that is on average more extreme than those outlets when visiting news sites through Facebook and 5% consume such extreme news when visiting news sites through other sources.

The analysis does not lead to conclusive results regarding changes in news consumption. Segregation online may have increased, but it probably did not change dramatically. How does this result coincide with increased segregation on social media? While I find that Facebook is more segregated than other online content, and while Facebook is typically the first or second most important source for online traffic, social media still accounts for a limited share of visits to news sites. For an average individual in the Comscore sample, 4% of news sites were visited through Facebook and in the extension subsample, which only includes Facebook users, the figure is 15%.²¹ Therefore, social media can be substantially more segregated than news consumed through other sources without dramatically changing overall segregation in online news consumption.

B Segregation Within Facebook

Why does news consumed through Facebook tend to be more extreme and segregated? Two mechanisms that could increase segregation are homophily in social networks (an “echo chamber” effect where one’s friends tend to recommend like-minded news sources) and the abundance of accessible, free media options allowing consumers to personalize their news feed. Panels 2 and 3 of Figure 4 show that the increased segregation is mostly associated with Facebook pages (the outlets participants subscribe to on Facebook) and not with Facebook friends (the social network). The isolation index is 0.14 when participants visit news sites not through Facebook, 0.18 when they visit sites through posts shared by Facebook friends, and 0.43 when they visit sites through posts shared by Facebook pages. Appendix Table A.7b shows that the results hold for the segregation measure as well.²²

This descriptive analysis cannot completely isolate each mechanism, nor rule out additional mechanisms. For example, posts by both Facebook friends and Facebook pages are also affected by Facebook’s algorithm, which is discussed in more detail in Section VI. Still, the analysis suggests that in order to understand segregation in social media, it is important to study the forces determining which pages appear in the Facebook feed and the effect of posts from these pages. Furthermore, posts shared by pages should not be ignored since approximately half of visits to news sites through Facebook in the extension subsample are through links shared by pages (row 10 in Appendix Table A.9b).

To conclude, in a 2019 survey, 83% of Americans stated that one-sided news is a very big or moderately big problem on social media.²³ This section provides evidence that this concern is warranted, as it shows that news accessed through Facebook is indeed more segregated and extreme than other online news. The next section estimates the causal effects of exposure to more and less segregated news using the random variation generated by the experiment.

IV Findings: Demand for News on Social Media

A Individuals Are Willing to Engage with Counter-Attitudinal News

Figure 6 displays the effects of the pro- and counter-attitudinal treatments on engagement with the potential pro- and counter-attitudinal outlets, respectively. To keep the results comparable across media outcomes, the figure is calculated for the participants who both installed the browser extension and provided permissions to access their posts for at least two weeks. Each row in the figure is estimated by regressing engagement with the four potential pro- or counter-attitudinal outlets in the two weeks following the intervention on the pro- or counter-attitudinal treatment. The control group is the reference group. Throughout the

²¹These estimates may underestimate Facebook usage since they rely on browsing activity on computers, while Facebook is more popular on mobile. For comparison, Parse.ly (2018) tracks pages viewed in thousands of sites and estimates that 16% of traffic related to Donald Trump in April-May 2018 was from social media and that Facebook is the largest external referral source for traffic in the law, government and politics category.

²²Figure 4 also provides a comparison of the isolation index in outlets individuals subscribe to on Facebook, posts they see in their feed, news sites they visit, and posts they share (Panel 1), and shows that isolation is highest among subscriptions.

²³Pew Research Center American Trends Panel Wave 51, July 2019.

analysis, I use linear regressions for ease of interpretation. Appendix Table A.10 shows that the effects on the feed, browsing behavior, and posts shared are qualitatively similar when running Poisson regressions.

The first panel of Figure 6 shows that the counter-attitudinal treatment increased the number of subscriptions to counter-attitudinal outlets by 1.42, compared to the control group. The effect is significant as the entire confidence interval is greater than zero. The increase is similar to the number of outlets participants immediately subscribed to in the intervention (1.51, not shown in the figure) since few participants unsubscribed from these outlets within two weeks.

Exposure to Posts in the Facebook Feed

The second panel of Figure 6 shows that in the two weeks following the intervention, participants in the pro- and counter-attitudinal treatments were exposed to 64 and 31 additional posts from the potential pro- and counter-attitudinal outlets, respectively. For comparison, control group participants were exposed to 266 posts from leading news outlets, and 2,335 posts in total, suggesting that the intervention affected news exposure but did not take over the participants' feeds.

Appendix Figure A.2 shows that the effect on exposure is driven mostly by organic posts published by pages and not by sponsored posts or posts shared by friends, meaning that participants were exposed to the content directly, without commentary from their social network. To test whether participants noticed the posts, they were asked in the endline survey how often they saw news from various outlets in their Facebook feed in the past week. Appendix Figures A.3 and A.4 show that participants reported seeing more news from the outlets they were offered and that participants in the counter-attitudinal treatment were less likely to say that opinions they see in their feed are aligned with their views. This implies that the effect on the feed was noticeable for at least two months, and confirms that the treatment affected the large subsample of participants who completed the endline survey and not only participants who installed the extension.

News Sites Visited

The third panel of Figure 6 shows that the counter-attitudinal treatment increased total visits to the websites of the counter-attitudinal outlets by 79%, an ITT effect of 1.34 visits over a baseline of 1.70 visits in the two weeks following the intervention. The pro-attitudinal treatment increased the number of visits to the websites of pro-attitudinal outlets by 21%, an ITT effect of 2.72 visits over a baseline of 13.21.

Appendix Figure A.5 separately estimates the effects of the intervention on the number of visits to the outlets' websites through a link appearing in the Facebook feed and on visits not directly associated with Facebook. While there is a strong and significant effect on visits through Facebook, there also seems to be an effect on other visits, albeit the latter result is not precisely estimated. It is possible that once participants read an article on the outlets' websites, they followed links to other articles as well. Alternatively, when participants became more familiar with the new outlets, they may have started visiting those outlets even without a Facebook referral. Appendix Figure A.6 shows that participants were more likely to click posts appearing higher in the feed. This could occur both because participants are more curious when they just start scrolling their feed and because Facebook's algorithm ranks posts according to expected interest. Interestingly, conditional on the order of posts, participants were as likely to visit a link from an outlet they subscribed to as a result of the intervention, compared to other news outlets.

Sharing Behavior

The fourth panel of Figure 6 shows that participants not only consumed news from counter-attitudinal outlets when they appeared in their feeds, they also shared the posts. To increase power, in Appendix Figure A.7, I analyze this effect using the entire access posts subsample and show that both treatments had a significant effect on the number of posts shared by these participants. The fact that participants chose to share the posts suggests that they considered the posts important, and implies that participants expanded the treatments to their social network.

Complementing previous studies focusing on Twitter (Halberstam and Knight, 2016), participants were much more likely to share pro-attitudinal posts. However, the relative effect on sharing counter-attitudinal posts compared to the control group (an increase of 105%) is stronger than the relative effect of the pro-attitudinal treatment (53%). Participants may have shared posts while commenting negatively on their

content. The second panel of Appendix Figure A.7 focuses on posts that were shared with no commentary by the participants and shows that even among these posts, the counter-attitudinal treatment had a significant effect on the number of posts shared.

B The Social Media Feed Strongly Affects Online News Consumption

The previous section demonstrated that individuals engage with the potential outlets when they appear in their feed, suggesting that news is often consumed incidentally when it becomes more accessible. This raises the question of whether individuals adjust the rest of their news consumption such that the slant of their news diet will not change. For example, individuals randomly offered the New York Times may start consuming more articles from the outlet's website, but consequently consume less news from the Boston Globe, which offers a similar perspective. To test whether the treatment affected the mean slant of all news participants engaged with, I focus on the conservative and liberal treatments since there are clear predictions on how these treatments would affect the slant.

Exposure to Posts on Facebook

The first panel of Figure 7 shows that when participants were randomly offered liberal or conservative outlets, their feed became substantially more liberal or conservative, respectively. The combined ITT effect of the liberal and conservative treatments equals 36% of the gap between the slant of the feed of liberals and conservatives in the control group, where slant is measured based on the leading news outlets dataset (participants who did not visit any news sites are excluded). The corresponding TOT effect is 47%. The change in slant provides a strong first stage, which is useful when analyzing the effect on political beliefs. It also allows me to test whether a change in the social media feed affects the slant of news sites visited or whether participants maintain a constant slant. The latter would suggest that participants re-optimize the sites they visit following an exogenous shock to their feed.

News Sites Visited

I find that individuals do *not* fully re-optimize their news consumption to keep the slant of the news sites they visit constant. The second panel of Figure 7 shows that the treatments had a strong and significant effect on the slant of news sites visited by the participants. The combined effects of the liberal and conservative treatments equal 14%-19% (ITT-TOT) of the difference in the slant of news sites visited by conservatives and liberals in the control group. Based on the Comscore panel, the TOT effect of the liberal treatment would shift the online news diet of an individual in Pennsylvania, a swing state, to a diet similar to an individual in New York, a blue state, and the TOT effect of the conservative treatment would shift the individual's news consumption to a news diet similar to an individual in South Carolina, a red state.²⁴ Appendix Table A.11 shows that the effect on slant is robust across various subsamples (e.g., when excluding participants who did not complete the endline survey).

By combining the exposure and browsing data, I find that when the compliers' news feed became one standard deviation more conservative, the slant of the news sites they visit became 0.31 standard deviations more conservative, and the slant of the subset of sites visited through Facebook became 0.71 standard deviations more conservative (both effects are significant at the 1% level). These estimates are calculated by instrumenting the slant of the posts observed in the Facebook feed with the treatment assignment. The regressions rely on the exclusion restriction that the treatments only affected the slant of sites visited through the slant of the Facebook feed. While the intervention is only expected to have an effect through the Facebook feed, the treatments could affect the feed in many ways. I am condensing the feed, a complicated object, to a scalar, the mean slant of news an individual was exposed to. This scalar is strongly affected by the treatment assignment and has intuitive economic meaning, but other changes in the feed, not captured in this measure, could affect the news sites visited. Since these calculations rely on stronger assumptions than the ITT and TOT estimates, they should be interpreted cautiously.

²⁴I determined the mean news consumption of each individual in Comscore's 2017 and 2018 panels based on visits to leading news outlets. Individuals who visited only one news site are excluded. The slant is then calculated at the state level for all panel members in the state.

To test for spillovers across news outlets, I calculate the effect of the treatments on the mean slant of all leading outlets excluding the eight potential outlets defined for each individual. Appendix Figure A.8 shows that the mean slant of news consumption is not strongly affected by the treatments when the potential outlets are excluded, implying that the experiment did not have large crowd-in or crowd-out effects.

Persistence

Is it possible that participants were initially curious about the new outlets they were offered but quickly stopped engaging with them. Figure 8 shows that the effect of the conservative treatment on news slant, compared to the liberal treatment, declines over the first six weeks after the intervention but mostly remains positive and significant. Appendix Figure A.9 repeats this analysis for the first twelve weeks after the intervention. While these results should be interpreted more cautiously since a substantial number of participants did not keep the extension installed or provide permissions to access posts over this longer time period, they suggest that the effects of the experiment declined but remained significant for at least twelve weeks.

The long-term effects also alleviate concerns that experimenter effects are driving the results in this section, as it is unlikely that participants remembered which posts appeared in their feed as a result of the intervention two months after the baseline survey, assumed that the experimenter expected them to persistently visit these websites, were constantly conscious that some of their browsing behavior could be observed, and were willing to spend time visiting news sites only to leave an impression on the experimenter. Furthermore, a survey question in the endline survey suggests that most participants did not remember which outlets they subscribed to and therefore their behavior or answers are unlikely to have been driven by experimenter effects.²⁵

C Discussion

This section shows that people are willing to substantially change their news consumption and engage with counter-attitudinal news on social media. Appendix C.2 analyzes the content of posts participants engaged with based on the words appearing in the posts and the article sections the posts linked to (e.g., Politics, Business, or Arts). I find that a large share of content tends to be political, even when the outlets the participants engaged with were counter-attitudinal.

How do these results coincide with the previous section, which shows that news consumed through social media, tends to be pro-attitudinal? If news is consumed incidentally on social media, and the Facebook feed tends to be pro-attitudinal, individuals are more likely to visit pro-attitudinal websites through social media but they will start visiting counter-attitudinal websites when they appear in their feed. Passive news consumption can also explain why Chen and Yang (2019) find that providing access to uncensored Internet does not lead to consumption of censored foreign news. As long as consumers are passive, providing access to new outlets may not be sufficient to affect news consumption because consumers will continue visiting the default outlets appearing in their bookmarks, search results, or social media feeds. My intervention may have affected news consumption because it increased the salience of specific outlets and decreased the search costs required to visit them by showing them on Facebook often.

This conclusion raises concerns regarding the power of social media companies in shaping news consumption habits. The effect of the social media feed on news consumption implies that any change to the feed, due to new subscriptions or a change in the algorithm, can drastically change one's news diet. Attempts to change the feed by suggesting new content happen all the time. They can stem from companies attempting

²⁵Participants were asked "In a previous survey, we may have asked if you are interested in 'liking' Facebook news pages. Did you like a page in the previous survey?" Only 40% of participants in the treatment arms stated that they remembered whether they liked a page and which pages they liked. Unfortunately, many participants did not understand this question and assumed it refers to a previous question in the endline survey. Therefore, I interpret this question as providing qualitative evidence that many participants did not remember which outlets they subscribe to and not for empirical analysis. The misunderstanding probably leads to an overestimation of the number of participants who remember which pages they liked as some respondents may have remembered the previous question in the endline survey but not the outlets offered in the baseline survey. Furthermore, even among the minority of participants who understood the question and stated that they remember which pages they liked, some did not state the correct outlets.

to maximize profits by increasing user engagement or originate from entities attempting to maximize political goals, such as political candidates purchasing ads or even foreign agents promoting Facebook pages to influence the American electorate.²⁶

V Findings: Opinions and Attitudes

A Social Media News Exposure Does Not Strongly Affect Political Opinions

The top panel of Figure 9 shows that the treatments did not affect the political opinions index. While the point estimate has the expected sign, the effect is minimal (0.005 standard deviations), precisely estimated, and not statistically significant. The upper bound for the combined liberal and conservative treatment effects, based on a 95% confidence interval, is only 0.8% of the difference in political opinions between liberals and conservatives in the control group. Appendix Figure A.10 shows that the effect on each component of the political opinions index is small, and I cannot reject a null effect for any of the components.

Why did the treatments not affect political opinions even though they dramatically affected the Facebook feed of participants? In other settings, studies on persuasion found a null effect that masked substantial heterogeneity (Baysan, 2019). Perhaps some participants were persuaded by the offered outlets, while for others, there was a backlash effect and opinions moved in the opposite direction of their treatment assignment. Appendix Figure A.11 estimates the effect of the interaction of ideology and treatment on the political opinions index and does not find evidence for a backlash effect. A second option is that the treatment did not affect political opinions since social media is still not a dominant news source, compared to television. This could explain why the results of this study differ from studies on Fox News (DellaVigna and Kaplan, 2007; Martin and Yurukoglu, 2017). Interestingly, I do not find evidence for heterogeneity based on whether participants reported getting most of their news on social media (see Appendix C.3). It is also possible that the null effects are explained by the fact that the intervention lasted for two months. However, the intervention lasted long enough to affect attitudes, as discussed in the next section.

The results differ from a recent study that found a backlash effect when exposing individuals to counter-attitudinal content on Twitter (Bail et al., 2018). Differences in the experiments' design can explain the differing results. Bail et al. (2018) expose individuals to bots retweeting counter-attitudinal *views*. Individuals plausibly become more upset when exposed to opposing opinion leaders, compared to counter-attitudinal news outlets. Bail et al. (2018) also provided monetary incentives to continuously follow the bots, asked participants to disable Twitter's timeline algorithm to ensure they viewed the tweets, and included weekly surveys to verify compliance. In my setting, participants could decide whether to comply with the treatment and engage with the content. Therefore, compliers with each treatment arm are different by design and this could affect the results. Social scientists have criticized the generalizability of forced exposure media experiments since the effects found may be concentrated among individuals who would not consume the content outside the experimental setting (Hovland, 1959; Bennett and Iyengar, 2008). For example, conservatives who get upset when visiting msnbc.com are less likely to consume content from MSNBC in my setting but may consume such content when encouraged to do so by the experimenter, and this type of consumption could drive the backlash effect.

B Exposure to Counter-Attitudinal News Decreases Affective Polarization

The bottom panel of Figure 9 shows that the counter-attitudinal treatment modestly decreased the affective polarization index compared to the pro-attitudinal treatment. The ITT and TOT effects are 0.03 and 0.06 standard deviations, respectively. This suggests that the concerns over more segregated news consumption are not misguided. When estimating the effect on each component of the index separately in Appendix

²⁶For example, many ads purchased by Russian organizations in their attempt to influence the 2016 election promoted Facebook pages. Congress has published the ads and they can be found here: <https://intelligence.house.gov/social-media-content/social-media-advertisements.htm>

Figure A.12, the effect is largest on whether participants find it difficult to see things from each party's point of view.

Appendix Tables A.12b, A.13, and A.14b show that the result is robust to excluding covariates, dropping each of the five components of the affective polarization measures from the index one at a time, and excluding participants who already subscribed to at least one of the primary outlets before the intervention. Appendix Table A.15b shows that an effect is detected when focusing on the subsample of participants who completed the endline survey and installed the extension. The effect is stronger among this group, which also had higher compliance rates. Appendix C.4 shows that the effect is similar when the regressions are reweighted to match populations means in ideology, party affiliation, gender, age, and the baseline feeling thermometer measure. Appendix C.5 estimates heterogeneous effects using causal forests (Wager and Athey, 2018) and shows that the predicted effect in the entire baseline sample is very similar to the effect among the endline survey subsample.

Comparing each treatment separately to the control group shows that most of the difference between the pro- and counter-attitudinal treatments stems from the counter-attitudinal treatment, perhaps because the relative effect of this treatment on engagement with the outlets was larger compared to baseline. In all specifications, the effect of the counter-attitudinal treatment is negative, statistically significant, and stronger than the effect of the pro-attitudinal treatment. However, this comparison suffers from differential attrition, due to lower attrition in the control group. Therefore, in Appendix Table A.12, I also calculate Lee bounds for the effects of each treatment (Lee, 2009). Due to the relatively small treatment effect, the bounds include a null effect. As an additional robustness test, I exclude control group participants who were recruited using the last email or ad inviting them to the endline survey (Behaghel et al., 2015). Without these participants, I compare the 45-46% of participants in each treatment arm who were "easiest" to recruit and attrition is similar across treatments. The results using this method are almost identical to the main specification.

I do not find evidence for substantial heterogeneity across most covariates I test for, including age, ideological leaning, baseline interest in news, and baseline exposure to counter-attitudinal news (Appendix C.3). One exception is that the treatment seemed to have a stronger effect on participants who were less polarized in baseline according to the feeling thermometer question. However, this effect is significant only at the 10% level and more research is required on heterogeneity.

In the rest of this section, I interpret the magnitudes of the effect using three approaches. First, I compare the effect of the intervention to benchmarks in the control group and outside the experiment. Second, I use the extension data to estimate the effect of a change in exposure to pro- and counter-attitudinal news on affective polarization. Third, I conduct two back-of-the-envelope calculations to estimate how affective polarization would have changed if Facebook had a more balanced feed. All the results are based on the effect of the offered outlets over two months and could be different with longer exposure or if different outlets were offered.

The ITT and TOT effects of the counter-attitudinal treatment decrease the difference between the feeling toward the participant's party and the opposing party by 0.58 and 0.96 degrees (on a 0-100 scale), respectively. For comparison, in the past 20 years, the feeling thermometer measure increased by 3.83-10.52 degrees. An additional point of comparison is a recent experiment by Allcott et al. (2020), who found that disconnecting from Facebook for one month in the fall of 2018 decreased the feeling thermometer measure by 2.09 degrees.²⁷ Hence, one way to interpret these results is that almost half of the depolarizing effect of disconnecting from Facebook can be achieved by replacing 1-4 subscriptions to pro-attitudinal outlets with subscriptions to counter-attitudinal outlets.

To estimate the effect of exposure to pro- or counter-attitudinal news on polarization, I focus on participants who installed the browser extension and completed the endline survey, i.e., the overlap between the extension and the endline subsamples. I use two summary measures for exposure to pro- and counter-attitudinal news: the share of counter-attitudinal news in the Facebook feed and the feed's congruence scale (both defined in Section II.D.1). I calculate these statistics based on all posts observed between the baseline and endline survey, for participants who observed at least two pro- or counter-attitudinal posts. I estimate the

²⁷Focusing on one measure decreases power. The effect I find on the feeling thermometer is statistically significant at the 10% level and the Allcott et al. (2020) benchmark is not statistically significant.

effect of each measure on affective polarization, and instrument the measure with the treatment assignment. Similar to the discussion in Section IV.B, the IV regressions rely on the exclusion restriction that the treatment only affects affective polarization through its effect on the measure analyzed.

I find that an increase of one standard deviation in the share of exposure to counter-attitudinal news decreases affective polarization by 0.13 standard deviations and an increase of one standard deviation in the congruence scale has a similar effect. The effects are significant at the 10% level as the sample size is smaller when focusing on participants who both installed the extension and completed the endline survey. One challenge in studying affective polarization based on non-experimental survey data (e.g., Garrett et al., 2014) is determining whether the correlation between news exposure and affective polarization is due to selection, i.e., individuals with more negative views of the opposing party select into more pro-attitudinal news exposure, or a causal effect, i.e., pro-attitudinal news makes people more polarized. Appendix Table A.16 shows that the effects of news exposure on affective polarization are approximately 26%-34% of the coefficients obtained using a cross-sectional regression among the control group, suggesting that the correlation is both due to a causal effect and selection.

I use the effect of the Facebook feed to estimate how affective polarization would have changed if individuals were exposed to more balanced news on Facebook. I find that if the feed had an equal share of pro- and counter-attitudinal news, the difference between the feelings toward one's party and the opposing party would decrease by 3.94 degrees. For this calculation, I estimate the effect of increasing the share of exposure to counter-attitudinal news by 33 percentage points, the difference between exposure in the control group and an exposure of 50%. The estimation does not rely on out-of-sample predictions as the share of counter-attitudinal news was greater than 50% for many participants in the counter-attitudinal treatment. Using a similar exercise, I find that if the congruence of the Facebook feed equaled zero, the difference between participants' feelings toward their one party and the opposing party would decrease by 3.43 degrees.

Perhaps a balanced news feed is not a realistic counterfactual because most individuals do not consume balanced news, regardless of social media. Therefore, in a second back-of-the-envelope calculation, I estimate how affective polarization would change if individuals were exposed in their Facebook feed to the same share of counter-attitudinal outlets, or the same congruence scale, as they encounter when visiting news sites not through Facebook. I find that the feeling thermometer outcome would decrease by 0.24-0.62 degrees. These calculations should be interpreted carefully since they do not take into account general equilibrium effects.²⁸ Nevertheless, they suggest that the Facebook feed may slightly amplify polarization.

B.1 Interpretation

Why did the treatments affect attitudes toward political parties but not political opinions? One possibility is that participants learned new facts about the world and these facts swayed their attitudes. Based on eight pre-registered survey questions, I test whether a change in participants' knowledge could explain the effect on polarization. In Appendix C.6, I do not find evidence for strong effects on knowledge.

Previous studies showed that Americans believe that members of the opposing party are more likely to hold extreme views than they actually do (Yudkin, Hawkins and Dixon, 2019), and therefore, attitudes may have changed because participants learned the opposing party is not as extreme as they thought.²⁹ I do not find evidence that the pro- and counter-attitudinal treatments had a significant effect on the distance between participants' baseline ideology and the perceived ideology of each party (Appendix Figure A.4).

Another option is that exposure to pro- and counter-attitudinal news affects attitudes due to increased negative coverage (Levendusky, 2013). This explanation predicts that pro-attitudinal outlets would increase

²⁸For example, it is likely that if Facebook drastically changed its feed, individuals would use other social media platforms instead. Some of this effect may be captured in the calculations since participants in the counter-attitudinal treatment used Facebook less often (as discussed in Section VI). However, with network effects, the decrease in Facebook use could be greater. The calculations also ignore the indirect effect of Facebook on news sites visited.

²⁹This theory is consistent with a study by Orr and Huber (2020) who find that negative attitudes toward individuals from the opposing party decrease when information is provided about their policy position.

negative attitudes toward the opposing party and counter-attitudinal outlets would affect consumers' attitudes toward their own party. This prediction is inconsistent with the data. I measure separately the effect of each treatment on attitudes toward each party and show in Appendix Table A.17 that the effect of the counter-attitudinal treatment on attitudes toward the opposing party is driving the results.

An alternative explanation, consistent with the data, is that participants exposed to counter-attitudinal news learned to rationalize the opinions of the opposing party. Intuitively, participants may have learned some of the opposing party's arguments and thus understood better why that party supports certain positions. This led to more positive attitudes but did not change political opinions as long as participants did not find these arguments particularly important. In Appendix D, I formalize this discussion using a model where political opinions are a weighted average of multiple beliefs and parties place different weights on beliefs.

There could be other explanations for the change in affective polarization.³⁰ The literature on affective polarization is new and more research is needed to pinpoint the precise mechanisms explaining how affective polarization evolves.

VI Findings: Exposure to Pro-Attitudinal News on Social Media

The previous section shows that exposure to pro-attitudinal news affects partisan hostility, therefore it is important to understand what influences the news individuals are exposed to on social media. This section decomposes the gap in exposure to posts shared by the pro- and counter-attitudinal outlets offered in the experiment into three main forces: participants are less likely to subscribe to counter-attitudinal news outlets; Facebook's algorithm supplies fewer posts from counter-attitudinal outlets, conditional on participants subscribing to them; and participants use Facebook less often when offered counter-attitudinal outlets. The decomposition exercise is based on the following framework:

$$E_{ij} = S_{ij}A_{ij}U_i$$

where E_{ij} , exposure, is the number of posts from outlet j that individual i was exposed to in her Facebook feed. Exposure is a product of whether individual i subscribed to outlet j (S_{ij}), the share of posts by the outlet among all posts the individual was exposed to, conditional on subscribing to the outlet (A_{ij}), and the total number of posts individual i was exposed to (U_i). I decompose the gap in exposure using the following formula:

$$\Delta E = \underbrace{S_{\Delta}A_C U_C}_{\text{Subscriptions}} + \underbrace{S_C A_{\Delta} U_C}_{\text{Algorithm}} + \underbrace{S_C A_C U_{\Delta}}_{\text{Usage}} + \underbrace{S_{\Delta} A_{\Delta} U_C + S_{\Delta} A_C U_{\Delta} + S_C A_{\Delta} U_{\Delta} + S_{\Delta} A_{\Delta} U_{\Delta}}_{\text{Combinations}} \quad (3)$$

where for each variable, $_C$ denotes the value for the counter-attitudinal treatment and $_{\Delta}$ denotes the difference between the pro- and counter-attitudinal treatments. *Subscriptions* is the additional counter-attitudinal posts participants assigned to the counter-attitudinal treatment would have been exposed to if they would have subscribed to the same number of outlets as participants assigned to the pro-attitudinal treatment. *Algorithm* is the additional posts these participants would have been exposed to if Facebook's algorithm would have supplied them with the same share of posts from counter-attitudinal outlets, as the share supplied when subscribing to pro-attitudinal outlets. *Usage* is the additional posts these participants would have been exposed to if they would have used Facebook as much as participants assigned to the pro-attitudinal treatment.

S_C and U_C are the mean number of new subscriptions and the total number of posts participants were exposed to, respectively, in the counter-attitudinal treatment. I estimate S_{Δ} and U_{Δ} by regressing the number

³⁰The counter-attitudinal treatment may have mitigated tribalism, which could have decreased affective polarization (Mason, 2015). Indeed, field experiments have found that strengthening partisan behavior can affect political behavior and beliefs (Gerber, Huber and Washington, 2010). I use party affiliation as a proxy for tribalism and find in Appendix Figure A.4 that the treatments did not significantly affect this proxy. However, the point estimate of the effect on Democratic Party affiliation has the predicted sign, and I cannot reject a small effect on affiliation with the Democratic Party.

of subscriptions and total exposure to posts on whether participants were assigned to the pro- or counter-attitudinal treatment. To estimate A_Δ and A_C , I pool the two groups of potential outlets for each participant such that each observation is a participant and either the group of pro-attitudinal outlets or the group of counter-attitudinal outlets. I then regress the share of posts that the participant was exposed to from a group of outlets (among all posts in the feed) on the full interaction of the number of new outlets the participant subscribed to and whether the group of outlets is pro-attitudinal. Since subscriptions are endogenous, they are instrumented with whether the group of outlets was randomly offered to the participant. The calculations are discussed in detail in Appendix C.7 along with alternative decompositions.

Figure 10 shows that the strongest force associated with participants' increased exposure to pro-attitudinal news is the algorithm. This demonstrates that even when individuals are willing to subscribe to outlets with a different point of view, Facebook's algorithm is less likely to show them content from those outlets (a phenomenon often described as a filter bubble). I also find evidence that participants prefer to subscribe to pro-attitudinal news outlets and that participants decrease their Facebook usage after they are offered to subscribe to counter-attitudinal outlets. The last effect is only significant at the 10% level and should be interpreted more cautiously. Still, it could explain why personalization is leading to segregation on social media. When consumers are exposed to more counter-attitudinal news, they may decrease their Facebook usage, and therefore, platforms have an incentive to filter counter-attitudinal news to maximize engagement. This result raises the question of whether the algorithm also personalizes content within an outlet, by showing conservatives more conservative posts shared by an outlet and liberals more liberal posts shared by the same outlet. In Appendix C.7.3, I find no evidence for within-outlet personalization.

Interestingly, even though I find that the algorithm seems to be filtering counter-attitudinal posts, Section III shows that the posts control group participants are exposed to in their feed are not more pro-attitudinal than the outlets they subscribe to on Facebook. One possible explanation for the differing results is that individuals probably subscribe to outlets as a response to non-random nudges. If nudges typically offer pro-attitudinal outlets, then users will subscribe to these outlets often and only users who are specifically interested in opposing content will subscribe to counter-attitudinal outlets. As a result, the algorithm may filter less counter-attitudinal content.³¹ The comparison to the control group descriptive statistics not only demonstrates why an experiment is necessary but also has policy implications. Adjusting the algorithm to offer more balanced news, conditional on subscription, would not make a big difference if individuals only subscribe to pro-attitudinal outlets. Therefore, to increase diversity in news exposure, nudges encouraging subscriptions to diverse outlets may also be required.

This section does *not* suggest that Facebook's algorithm intentionally increases segregation by ranking posts according to whether they match the user's beliefs, or that the interaction of the slant of an outlet and ideology of a user has a causal effect on a post's ranking. Platforms rank posts based on many signals that can be correlated with whether an outlet is counter-attitudinal, including the consumer's past engagement with the outlet, her social network, and possibly other pages she subscribes to. In other words, the effect of the algorithm also captures the behavior and perceived interests of the user. Indeed, Appendix C.7.2 shows that the effect of the algorithm slightly increases over time, suggesting that engagement with content plays a role in the ranking of posts.

Personalization of news exposure is still an important departure from how news was supplied in the past. Until recently, the engagement of an individual with news, e.g., the articles she read in the newspaper or the cable channels she chose to watch, did not affect her supply of news.

While I focus on Facebook, this section's conclusions likely apply to other platforms personalizing content as well. For example, since at least 2016, Twitter has been ranking tweets according to how interesting they would be for a user, based on factors such as the user's past relationship with the author. Therefore, it is plausible that tweets from pro-attitudinal accounts will receive a higher ranking. Furthermore, major news outlets have also started to personalize their websites and the articles they suggest to their customers.

³¹The control group participants subscribing to pro- and counter-attitudinal outlets are substantially different from each other. For example, among the twenty most popular liberal and conservatives pages, there is a difference of 0.32 standard deviations in the absolute value of ideology of subscribers to at least one pro- and counter-attitudinal outlet. Moreover, subscriptions to counter-attitudinal outlets occur several months later than subscriptions to pro-attitudinal outlets, and posts from more recent subscriptions are probably more likely to appear in the feed. The experiment assures that all subscriptions occur at the same time and due to a random offer.

VII Conclusions

Consumption of news through social media is increasing, but the effect of social media on public opinion remains controversial. I show that news consumption on social media is an important phenomenon because consumers are exposed to different news on social media, individuals incidentally consume news when it becomes accessible in their feed, and exposure to news on social media affects attitudes.

This paper supports a more nuanced view regarding the effect of media on public opinion. On the one hand, I show that exposure to pro-attitudinal news increases affective polarization compared to counter-attitudinal news. This result provides a mechanism complementing other important studies finding that social media can increase polarization and raises concern since affective polarization may decrease trust in government and the accountability of elected officials. On the other hand, individuals are not easily persuaded by the political leaning of their news exposure. The results of the experiment are in line with the long term increase in affective polarization, without an equivalent change in political opinions (Mason, 2015). This suggests that a more segregated news environment may partially explain the increase in affective polarization over the past decades.

Methodologically, this paper has several limitations. First, I only observe online news consumption. While I show that the intervention did not have substantial spillovers across online outlets, to precisely measure total news consumption, future studies would need to collect consumption data from other mediums, such as television, as well. Furthermore, I collect data on browsing behavior and the Facebook feed on a computer, but a growing share of news is consumed through smartphones. Second, while I argue that due to the organic nature of the intervention, it is unlikely that experimenter effects play a major role in this study, I cannot rule out that the perceived expectations of the experimenter affected the results. Third, the endline survey suffers from high attrition. I use several methods to alleviate this concern, but attrition could still affect the survey outcomes. Fourth, the study does not generate random variation in exposure to moderate outlets and therefore, cannot speak to their effects. Fifth, while the experiment has high external validity when it comes to analyzing partisan outlets on Facebook in 2018, the result may not hold for other periods. For example, Trump's presidency is exceptional in the stability of the president's approval ratings. If other opinions were relatively stable throughout the period as well, the null effect on political opinions could be explained by the period when the survey took place. Finally, I estimate all effects over several weeks or months, and the results may be different in the long-term.

This study has important policy implications. I demonstrate that Facebook's algorithm limits exposure to counter-attitudinal news. Automated personalization of news content may have stronger impacts in the future, due to growth in online news consumption and advances in machine learning algorithms customizing news exposure. However, I also find that individuals are willing to engage with counter attitudinal news. Therefore, even though social media platforms are associated with pro-attitudinal content, they can expose individuals to more perspectives. Suggestions include making algorithms more transparent, nudging users to diversify their feed, and modifying algorithms to encourage serendipitous encounters (Pariser, 2011; Sunstein, 2017). The experiment described in this paper essentially measures the effect of one such intervention and shows that a simple scalable nudge can effectively diversify news exposure and decrease polarization.

While social media algorithms may increase affective polarization through their effect on news consumption, platforms also have the potential to mitigate these effects.

References

- Allcott, Hunt, Luca Braghieri, Sarah Eichmeyer, and Matthew Gentzkow. 2020. "The Welfare Effects of Social Media." *American Economic Review*, 110(3): 629–676.
- Bail, Christopher, Lisa Argyle, Taylor Brown, John Bumpus, Haohan Chen, M.B. Hunzaker, Jaemin Lee, Marcus Mann, Friedolin Merhout, and Alexander Volfovsky. 2018. "Exposure to Opposing Views can Increase Political Polarization: Evidence from a Large-Scale Field Experiment on Social Media." *Proceedings of the National Academy of Sciences of the United States of America*, 115(37): 9216–9221.
- Bakshy, Eytan, Solomon Messing, and Lada A Adamic. 2015. "Exposure to Ideologically Diverse News and Opinion on Facebook." *Science*, 348(6239): 1130–1132.
- Baysan, Ceren. 2019. "The Polarizing Effects of Persuasive Communication: Experimental Evidence from Turkey."
- Behaghel, Luc, Bruno Crépon, Marc Gurgand, and Thomas Le Barbanchon. 2015. "Please Call Again: Correcting Nonresponse Bias in Treatment Effect Models." *Review of Economics and Statistics*, 97(5): 1070–1080.
- Bennett, W. Lance, and Shanto Iyengar. 2008. "A New Era of Minimal Effects? The Changing Foundations of Political Communication." *Journal of Communication*, 58(4): 707–731.
- Berinsky, Adam J., Gregory A. Huber, and Gabriel S. Lenz. 2012. "Evaluating Online Labor Markets for Experimental Research: Amazon.com's Mechanical Turk." *Political Analysis*, 20(3): 351–368.
- Boxell, Levi, Matthew Gentzkow, and Jesse M. Shapiro. 2018. "Greater Internet Use is Not Associated with Faster Growth in Political Polarization among US Demographic Groups." *Proceedings of the National Academy of Sciences of the United States of America*, 115(3): 10612–10617.
- Bursztyn, Leonardo, Georgy Egorov, Ruben Enikolopov, and Maria Petrova. 2019. "Social Media and Xenophobia: Evidence from Russia." *Working Paper*.
- Chen, Yuyu, and David Y. Yang. 2019. "The Impact of Media Censorship: 1984 Or Brave New World?" *American Economic Review*, 109(6): 2294–2332.
- Chiang, Chun Fang, and Brian Knight. 2011. "Media Bias and Influence: Evidence from Newspaper Endorsements." *Review of Economic Studies*, 78(3): 795–820.
- Comscore. 2018. *Web Behavior Database Panel 2007, 2008, 2017, 2018*. Wharton Research Data Services, University of Pennsylvania.
- Coppock, Alexander, Emily Ekins, and David Kirby. 2018. "The Long-lasting Effects of Newspaper Op-Eds on Public Opinion." *Quarterly Journal of Political Science*, 13(1): 59–87.
- Davis, Mark H. 1980. "A Multidimensional Approach to Individual Differences in Empathy." *JSAS Catalog of Selected Documents in Psychology*, 10.
- DellaVigna, Stefano, and Ethan Kaplan. 2007. "The Fox News Effect: Media Bias and Voting." *The Quarterly Journal of Economics*, 122(3): 1187–1234.
- Durante, Ruben, Paolo Pinotti, and Andrea Tesei. 2019. "The Political Legacy of Entertainment TV." *American Economic Review*, 109(7): 2497–2530.
- Enikolopov, Ruben, Alexey Makarin, and Maria Petrova. 2020. "Social Media and Protest Participation: Evidence from Russia." *Econometrica*, 88(4): 1479–1514.
- Flaxman, Seth R, Goel Sharad, and Justin M Rao. 2016. "Filter Bubbles, Echo Chambers, and Online News Consumption." *Public Opinion Quarterly*, 80: 298–320.

- Garrett, R. Kelly, Shira Dvir Gvirsman, Benjamin K. Johnson, Yariv Tsfati, Rachel Neo, and Aysenur Dal.** 2014. "Implications of Pro- and Counterattitudinal Information Exposure for Affective Polarization." *Human Communication Research*, 40(3): 309–332.
- Gentzkow, Matthew, and Jesse M. Shapiro.** 2010. "What Drives Media Slant? Evidence From U.S. Daily Newspapers." *Econometrica*, 78(1): 35–71.
- Gentzkow, Matthew, and Jesse M. Shapiro.** 2011. "Ideological Segregation Online and Offline." *Quarterly Journal of Economics*, 126(4): 1799–1839.
- Gentzkow, Matthew, Jesse M. Shapiro, and Daniel F. Stone.** 2015. "Media Bias in the Marketplace: Theory." In *Handbook of Media Economics*, 1B. Vol. 1, 623–645. Elsevier B.V.
- Gentzkow, Matthew, Jesse M. Shapiro, and Michael Sinkinson.** 2011. "The Effect of Newspaper Entry and Exit on Electoral Politics." *American Economic Review*, 101: 2980–3018.
- Gerber, Alan S., Dean Karlan, and Daniel Bergan.** 2009. "Does the Media Matter? A Field Experiment Measuring the Effect of Newspapers on Voting Behavior and Political Opinions." *American Economic Journal: Applied Economics*, 1(2): 35–52.
- Gerber, Alan S., Gregory A. Huber, and Ebonya Washington.** 2010. "Party Affiliation, Partisanship, and Political Beliefs: A Field Experiment." *American Political Science Review*, 104(4): 720–744.
- Gosling, Samuel D., Peter J. Rentfrow, and William B. Swann.** 2003. "A Very Brief Measure of the Big-Five Personality Domains." *Journal of Research in Personality*, 37(6): 504–528.
- Guess, Andrew.** 2020. "(Almost) Everything in Moderation: New Evidence on Americans' Online Media Diets." *American Journal of Political Science*, , (Forthcoming).
- Guess, Andrew, Brendan Nyhan, and Jason Reifler.** 2017. "You're Fake News" Findings from the Poynter Media Trust Survey. The Poynter Ethics Summit.
- Guess, Andrew, Brendan Nyhan, Benjamin Lyons, and Jason Reifler.** 2018. *Avoiding the Echo Chamber about Echo Chambers*. Knight Foundation.
- Halberstam, Yosh, and Brian Knight.** 2016. "Homophily, Group Size, and the Diffusion of Political Information in Social Networks: Evidence from Twitter." *Journal of Public Economics*, 143: 73–88.
- Hortacsu, Ali, Matthijs R Wildenbeest, and Bar De Los Santos.** 2012. "Testing Models of Consumer Search using Data on Web Browsing and Purchasing Behavior." *American Economic Review*, 102: 2955–2980.
- Hovland, Carl I.** 1959. "Reconciling Conflicting Results Derived from Experimental and Survey Studies of Attitude Change." *American Psychologist*, 14(1): 8–17.
- Iyengar, Shanto, and Masha Krupenkin.** 2018. "The Strengthening of Partisan Affect." *Political Psychology*, 39: 201–218.
- Iyengar, Shanto, Yphtach Lelkes, Matthew Levendusky, Neil Malhotra, and Sean J. Westwood.** 2019. "The Origins and Consequences of Affective Polarization in the United States." *Annual Review of Political Science*, 22(1): 129–146.
- Jo, Donghee.** 2018. "Better the Devil You Know: An Online Field Experiment on News Consumption."
- Kennedy, Patrick J., and Andrea Prat.** 2019. "Where Do People Get Their News." *Economics Policy Journal*, 5-27.
- Lee, David S.** 2009. "Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects." *Review of Economic Studies*, 76(3): 1071–1102.
- Lelkes, Yphtach, Gaurav Sood, and Shanto Iyengar.** 2015. "The Hostile Audience: The Effect of Access to Broadband Internet on Partisan Affect." *American Journal of Political Science*, 61(1): 5–20.

- Levendusky, Matthew.** 2013. "Partisan Media Exposure and Attitudes Toward the Opposition." *Political Communication*, 30(4): 565–581.
- Martin, Gregory J., and Ali Yurukoglu.** 2017. "Bias in Cable News: Persuasion and Polarization." *American Economic Review*, 107(9): 2565–2599.
- Mason, Lilliana.** 2015. "'I Disrespectfully Agree': The Differential Effects of Partisan Sorting on Social and Issue Polarization." *American Journal of Political Science*, 59(1): 128–145.
- Mosquera, Roberto, Mofioluwasademi Odunowo, and Trent Mcnamara.** 2019. "The Economic Effects of Facebook." *Experimental Economics*, 1–28.
- Müller, Karsten, and Carlo Schwarz.** 2019. "From Hashtag to Hate Crime: Twitter and Anti-Minority Sentiment."
- Okuyama, Yoko.** 2019. "Toward Better Informed Decision-Making: the Impacts of a Mass Media Campaign on Women's Outcomes in Occupied Japan."
- Orr, Lilla V, and Gregory A Huber.** 2020. "The Policy Basis of Measured Partisan Animosity in the United States." *American Journal of Political Science*, 64(3): 569–586.
- Pariser, Eli.** 2011. *The Filter Bubble*. The Penguin Press.
- Parse.ly.** 2018. "The Authority Report: 2018 Traffic Sources by Content Categories and Topics."
- Peterson, Erik, Goes Shared, and Shanto Iyengar.** 2019. "Partisan Selective Exposure in Online News Consumption: Evidence from the 2016 Presidential Campaign." *Political Science Research and Methods*, 1–17.
- PEW.** 2015. *Millennials & Political News: Social Media - the Local TV for the Next Generation*. Pew Research Center.
- Rand, David G., Alexander Peysakhovich, Gordon T. Kraft-Todd, George E. Newman, Owen Wurzbacher, Martin A. Nowak, and Joshua D. Greene.** 2014. "Social Heuristics Shape Intuitive Cooperation." *Nature Communications*, 5: 1–12.
- Reuters Institute.** 2019. *Digital News Report 2019*. University of Oxford.
- Strömberg, David.** 2015. "Media and Politics." *Annual Review of Economics*, 7(1): 173–205.
- Sunstein, Cass.** 2017. *#Republic: Divided Democracy in the Age of Social Media*. Princeton University Press.
- Wager, Stefan, and Susan Athey.** 2018. "Estimation and Inference of Heterogeneous Treatment Effects using Random Forests." *Journal of the American Statistical Association*, 113(523): 1228–1242.
- Williamson, Debra Aho.** 2018. *US Time Spent with Social Media 2019*. eMarketer.
- Yudkin, Daniel, Stephen Hawkins, and Tim Dixon.** 2019. *The Perception Gap: How False Impressions are Pulling Americans Apart*. More In Common.
- Zhuravskaya, Ekaterina, Maria Petrova, and Ruben Enikolopov.** 2020. "Political Effects of the Internet and Social Media." *Annual Review of Economics*.

Figure 1: Experimental Design

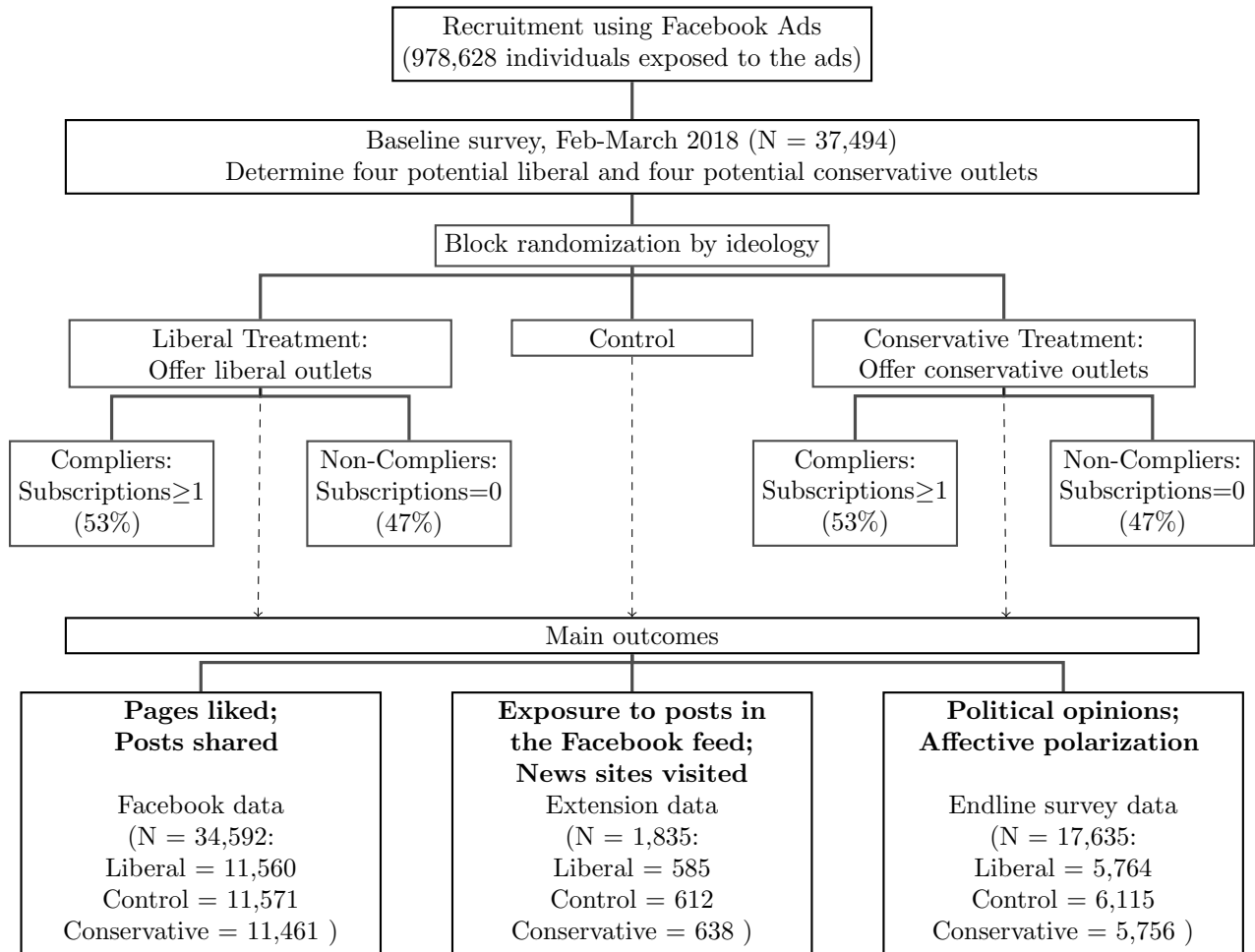
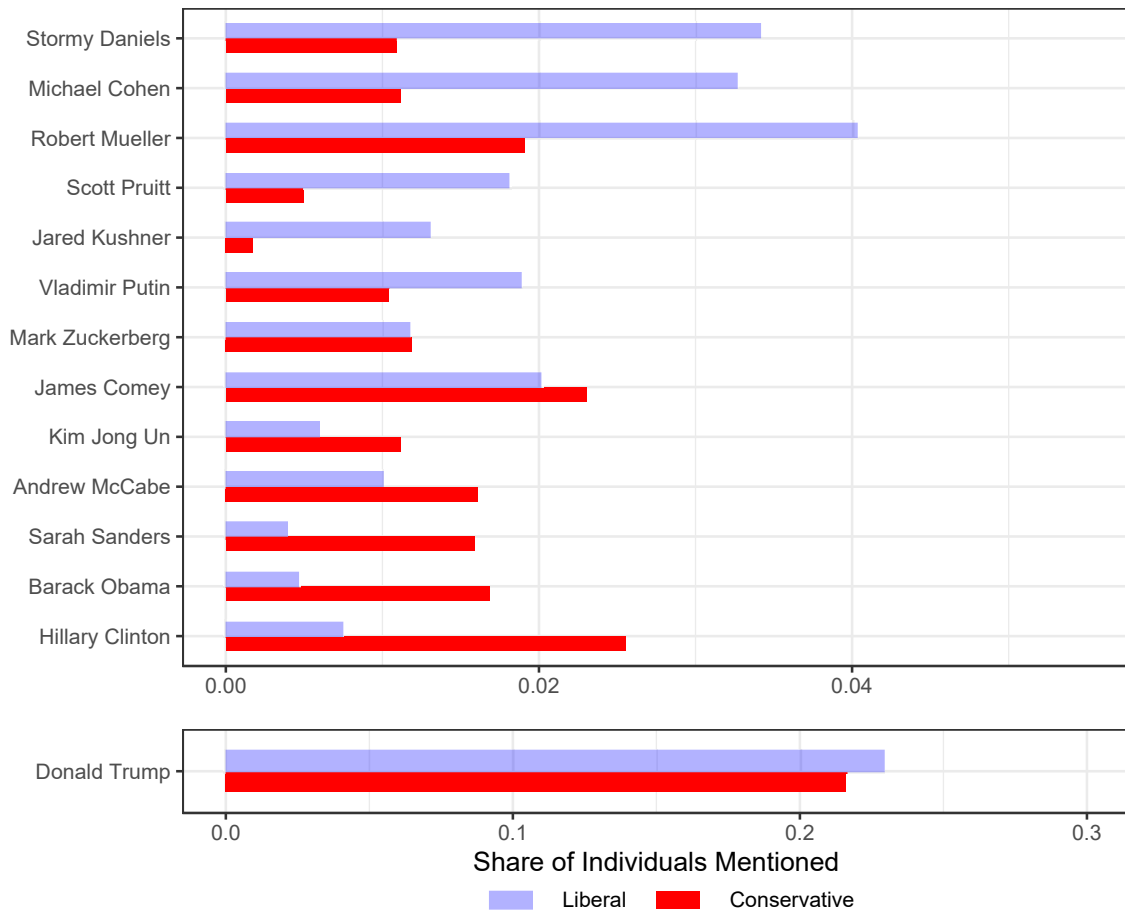
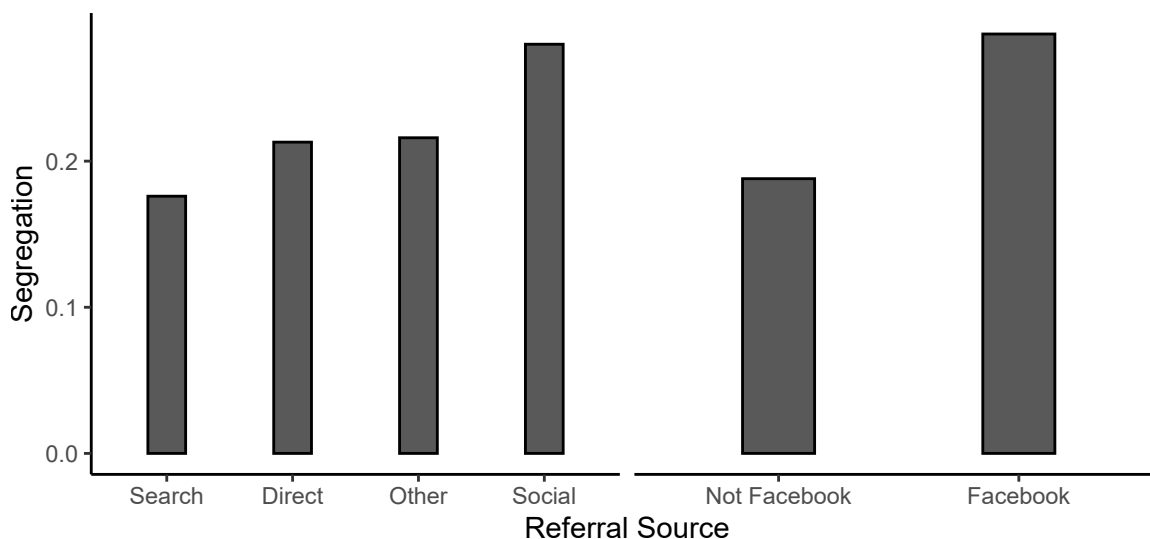


Figure 2: Figures Mentioned in Posts Shared by Outlets During the Study Period



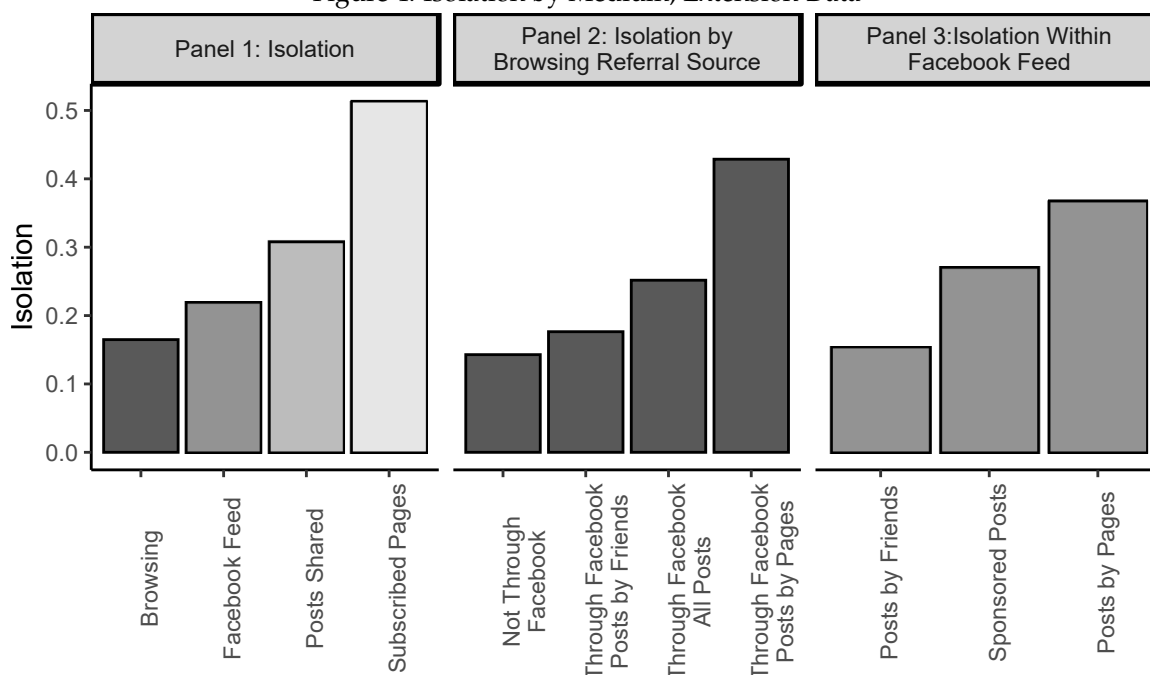
This figure shows the men and women mentioned most often in posts shared by the eight primary outlets and two main alternative outlets between February 28 and April 25, 2018, the median dates the baseline survey and endline survey were taken. Approximately 32% of posts with text mentioned a name. The x-axis is the share of times an individual was mentioned in a post by one of the liberal outlets (top bars) and by one of the conservative outlets (bottom bars), of all mentions of individuals. To fit all the figures on the same scale, the x-axis is broken for Donald Trump, who is by far the most dominant person mentioned. The figures were identified using the Spacy Natural Language Processing algorithm and post-processing names (e.g., removing possessive 's). Names that appear in only one outlet are excluded. If only a last name is mentioned, it is associated with the dominant first and last name combination when such a combination exists. To simplify the graph, the names 'Trump' and 'Donald Trump' are determined to be the same individual, even though 'Trump' could refer to other members of President Trump's family.

Figure 3: Segregation in News Sites Visited by Referral Source, Comscore Data



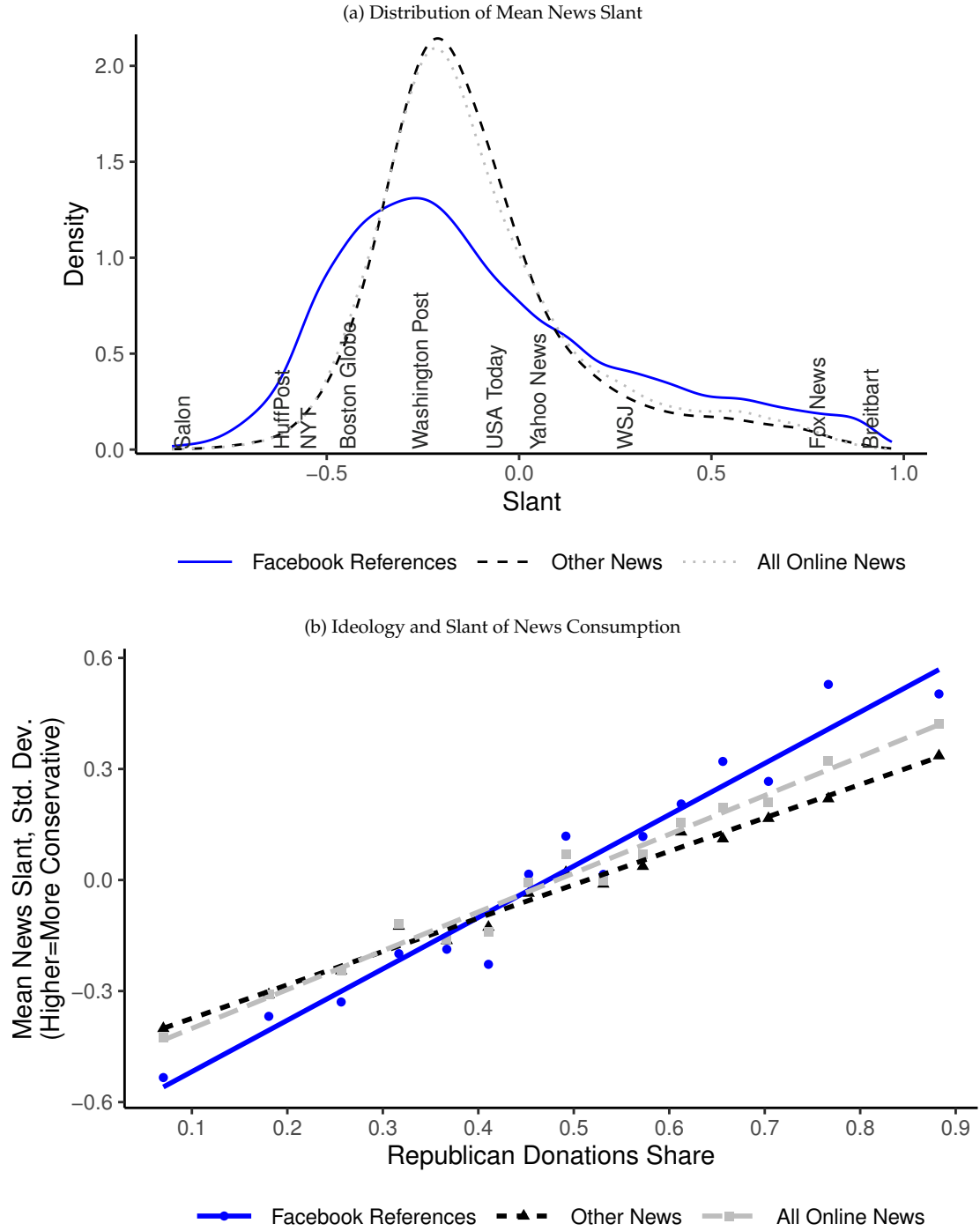
This figure displays the segregation in news sites visited by referral source. The definition of the segregation measure is discussed in Section III. Appendix A.5 defines the websites composing each channel. 2017-2018 Comscore data.

Figure 4: Isolation by Medium, Extension Data



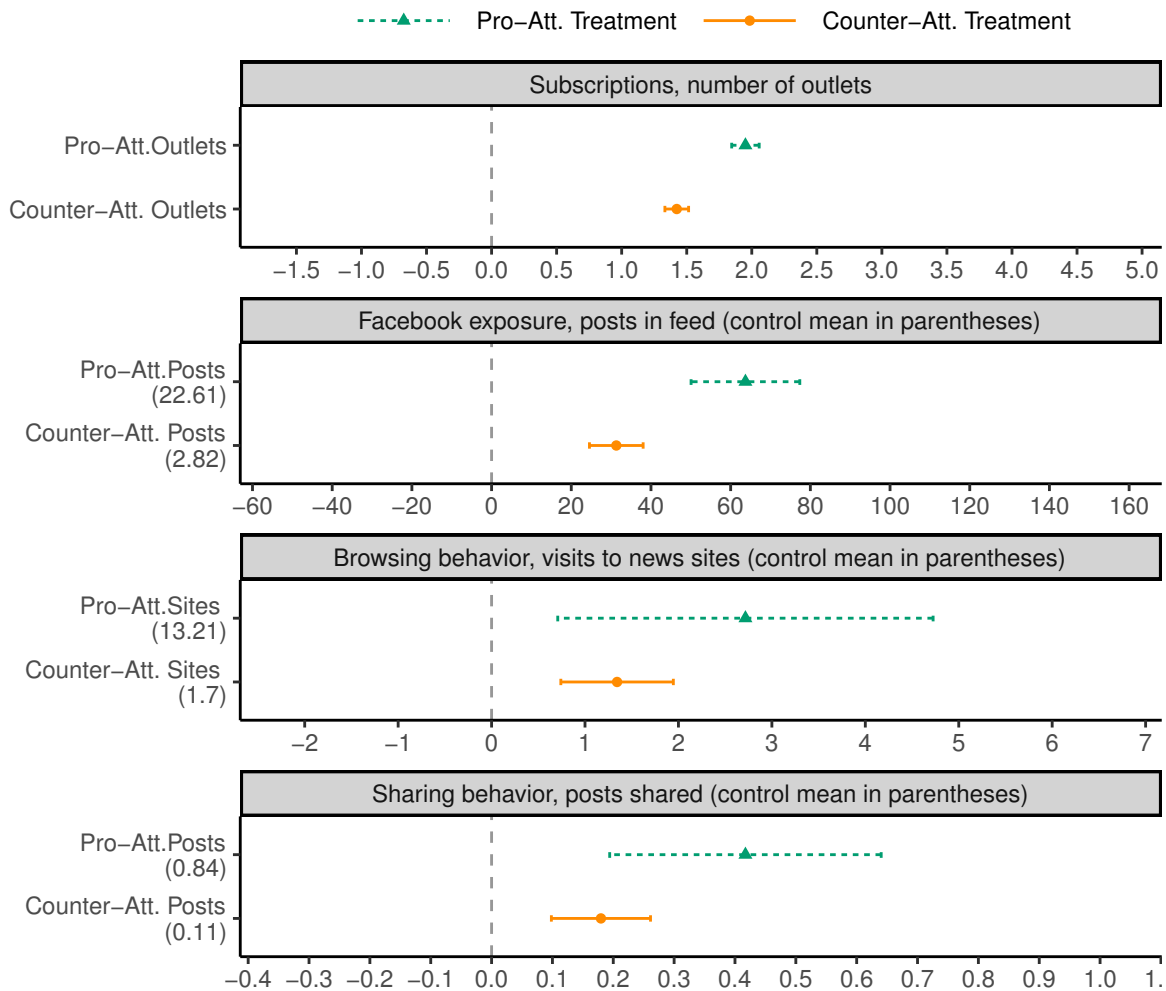
This figure displays the isolation of news participants engaged with. A higher value means liberals and conservatives were more likely to engage with different news outlets. Panel 1 shows the isolation measure for news sites participants visited, posts that appeared in their feed, posts they shared, and news outlets they subscribed to on Facebook. Panel 2 compares isolation values for news sites visited through different sources. Panel 3 compares different types of posts in the Facebook feed. The figure analyzes data from control group participants in the first eight weeks after the extension was installed.

Figure 5: News Consumption in the Comscore Panel



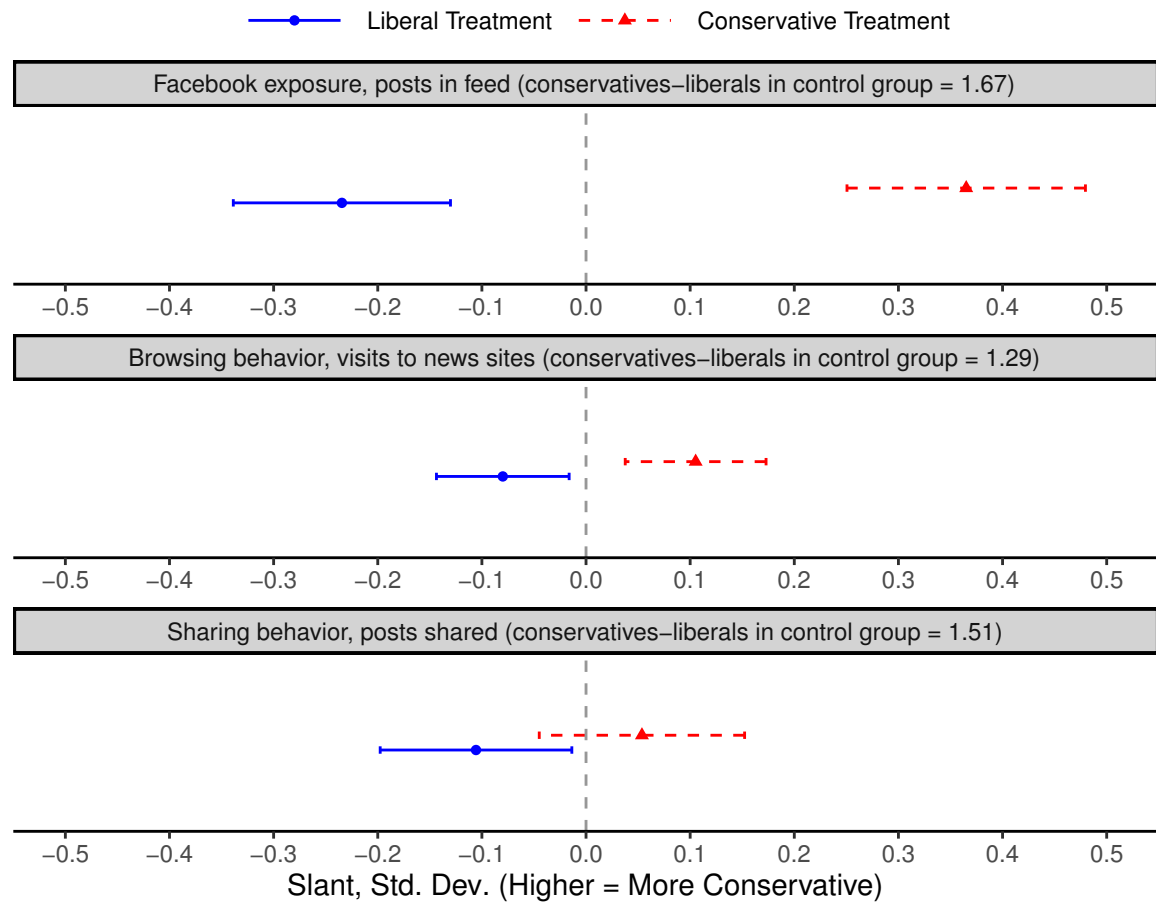
Sub-figure (a) presents the distribution of the mean slant of news sites visited (smoothing bandwidth = 0.05). Major news outlets are added to the x-axis for reference. The slant of each domain is based on Bakshy, Messing and Adamic (2015). A visit is referred from Facebook if the referring domain is “facebook.com.” Sub-figure (b) presents a binned scatter plot. The x-axis is the share of Republican donations in a zip code based on FEC donation data for the 2016 and 2018 election cycles and the y-axis is the mean slant of news sites visited. The sample for both figures includes individuals in the 2017 and 2018 Comscore Web Behavior Database Panel who visited news sites multiple times through Facebook and through other sources.

Figure 6: Effects of the Pro- and Counter-Attitudinal Treatments on Subscriptions, News Exposure, News Sites Visited and Sharing Behavior, Two Weeks Following the Intervention



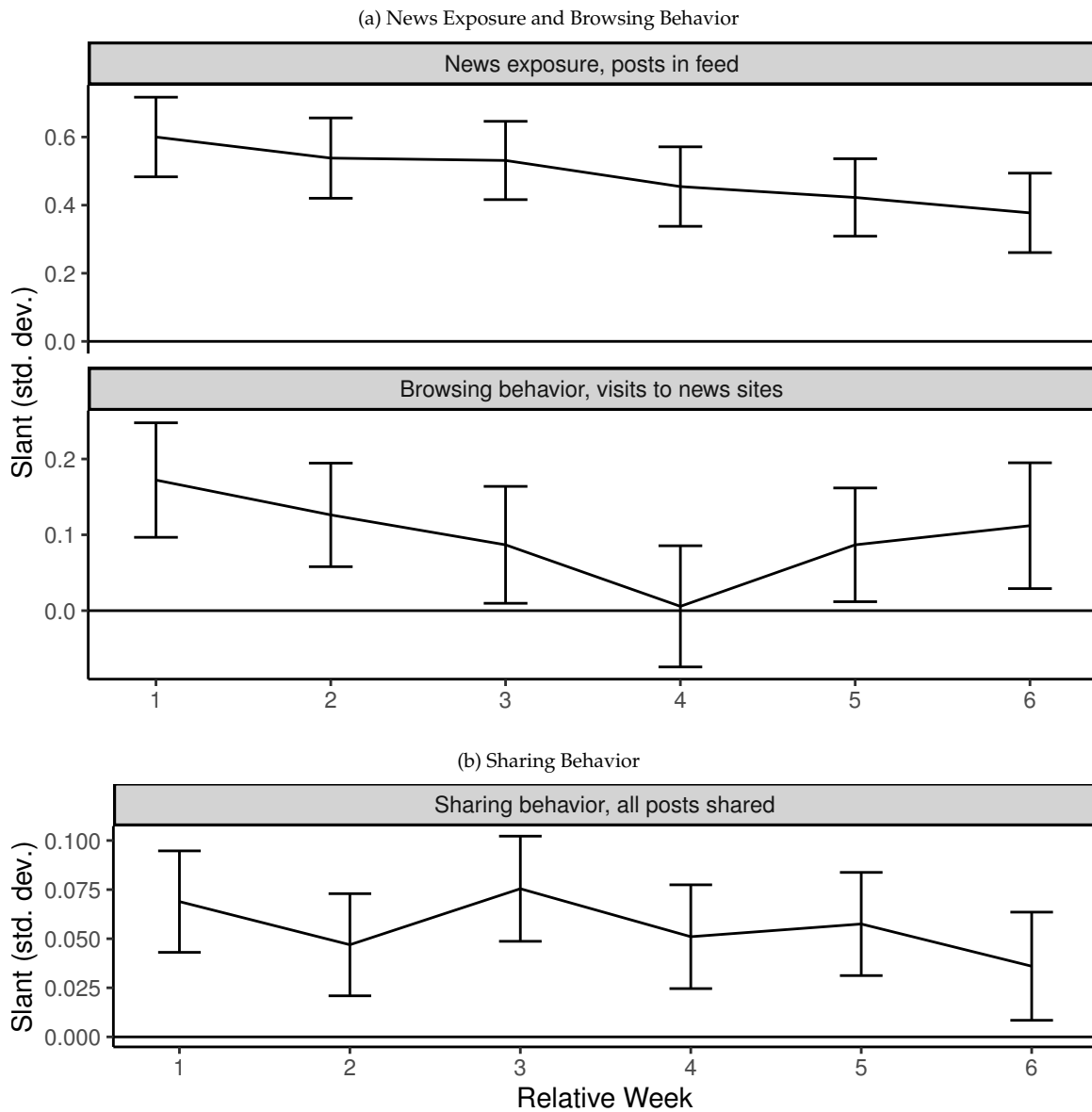
This figure shows the effect of the treatments on engagement with the offered outlets in the two weeks following the intervention. The dependent variable is engagement with either the four potential pro-attitudinal outlets or the four potential counter-attitudinal outlets and the independent variable is the treatment. Each panel presents the effect of a separate outcome. For example, in the third panel, the triangle and dashed line present the point estimate and the confidence interval of the effect of the pro-attitudinal treatment on visits to the websites of the potential pro-attitudinal outlets, compared to the control group. The regressions control for the outcome measure in baseline if it exists. The sample includes 1,648 participants with a liberal or conservative ideological leaning who installed the extension and provided permissions to access their posts for at least two weeks. Error bars reflect 90 percent confidence intervals.

Figure 7: Effect of the Treatments on News Slant



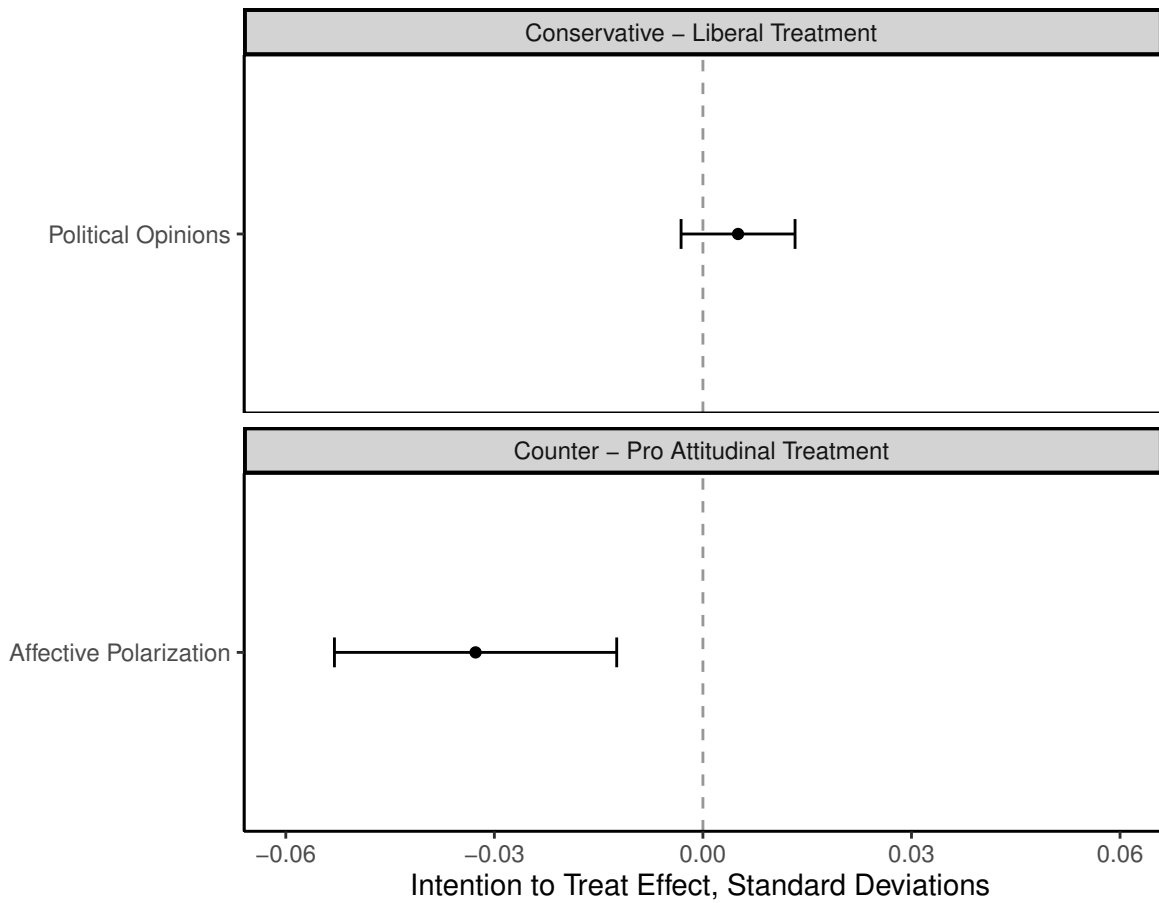
This figure shows the effect of the liberal and conservative treatments on the mean slant, in standard deviations, of all news individuals engaged with. In each panel, the dependent variable is the mean slant of outlets and the independent variable is the treatment. The regressions control for the outcome in baseline, if it exists. The sample includes participants who installed the extension and provided permissions to access their posts for at least two weeks following the intervention. Error bars reflect 90 percent confidence intervals.

Figure 8: Effects of the Conservative Treatment on Mean Slant by Week, Compared to the Liberal Treatment



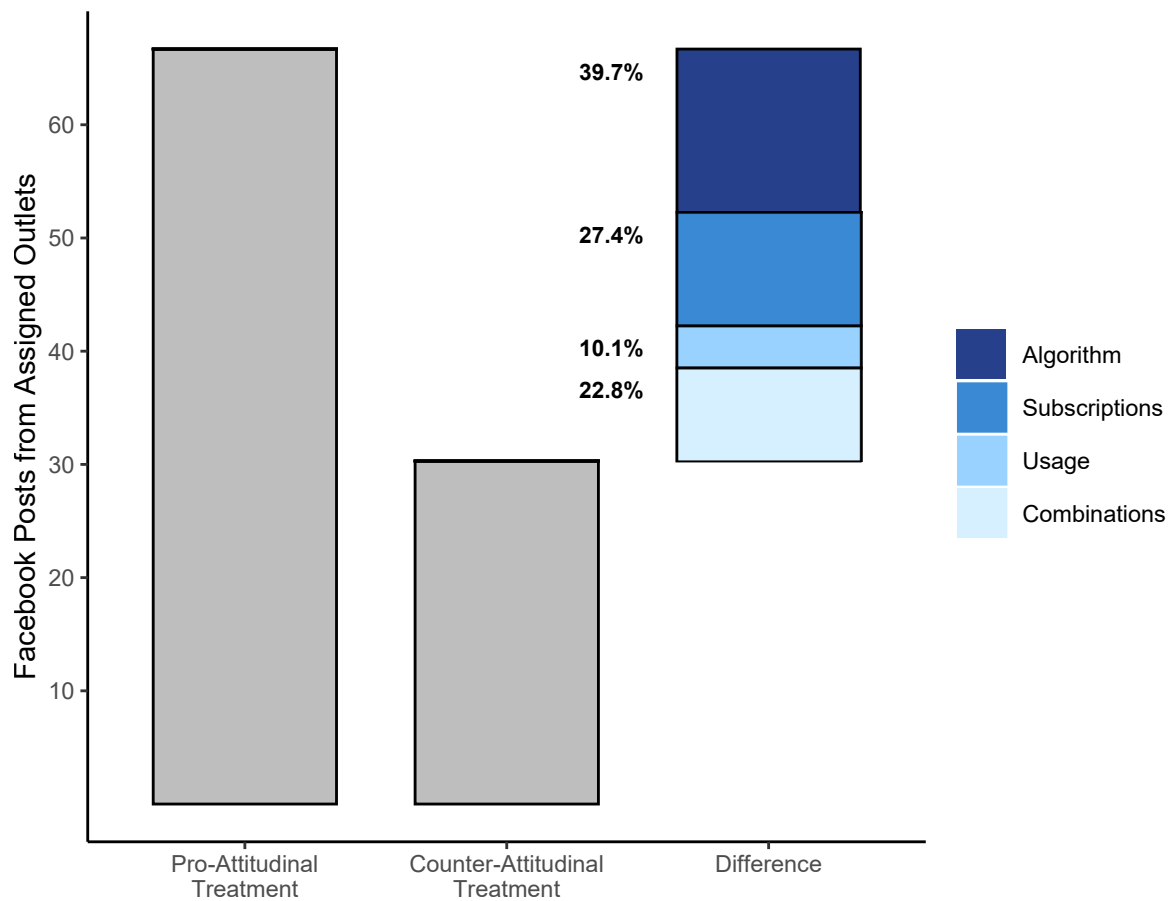
These figures show the difference between the effect of the liberal and conservative treatments on the mean slant of news engagement over time. Each panel presents a series of regressions, where the dependent variable is the slant of outlets in a specific week. The regressions control for the outcome in baseline when it exists. In the x-axis, relative week 1 is a full week immediately following the intervention. Sub-figure (a) is based on 1,596 participants who kept the extension installed for at least six weeks following the intervention. Sub-figure (b) is based on 29,131 participants who provided access to posts they shared for at least six weeks. Error bars reflect 90 percent confidence intervals.

Figure 9: Effect of the Treatments on Political Opinions and Affective Polarization



This figure shows the effect of the treatments on the primary endline survey outcomes. The first panel shows the effect of the conservative treatment on the political opinions index, compared to the liberal treatment. A higher value is associated with a more conservative outcome. The second panel shows the effect of the counter-attitudinal treatment on the affective polarization index, compared to the pro-attitudinal treatment. A higher value is associated with a more polarized outcome. The indices are described in Section II.D.2 and the regressions specifications are detailed in Section II.E. The panels are based on 17,635 participants who took the endline survey. Error bars reflect 90 percent confidence intervals.

Figure 10: Decomposing the Gap Between Exposure to Posts from the Offered Pro-Attitudinal and Counter-Attitudinal Outlets



This figure decomposes the gap between the number of posts participants were exposed to from the offered pro- and counter-attitudinal outlets. The y-axis is the number of posts seen in the feed in the two weeks following the intervention and the x-axis is the treatment arm. *Algorithm* describes the gap explained by Facebook's tendency to show participants a greater share of posts from pro-attitudinal outlets (among all posts in the feed) conditional on subscriptions. *Subscriptions* describes the gap explained by participants' tendency to subscribe to more offered outlets in the pro-attitudinal treatment. *Usage* describes the gap explained by participants' tendency to view fewer posts on Facebook (use Facebook less often) in the counter-attitudinal treatment. *Combinations* describe interactions between these expressions. Data is based on 1,059 participants in the pro- and counter-attitudinal treatments for which posts in the Facebook feed could be observed in the two weeks following the intervention and at least one post is observed. The calculations appear in Appendix C.7.

Table 1: Samples, Data Sources and Outcomes

Sample	Data Sources	Number of Participants and Retention	Main Outcomes
Baseline sample	Baseline survey; Facebook data on participants' subscriptions to outlets	37,494 (all participants)	Subscriptions to outlets in the intervention (compliance)
Access posts subsample	Facebook data for participants who provided permissions to access their posts and subscriptions for at least two weeks	34,592 (94% of participants who provided permissions in baseline)	Subscription to outlets over time; posts shared by participants
Extension subsample	Browser data for participants who installed the extension for at least two weeks	1,835 (81% of participants who installed the extension in baseline)	Exposure to posts in the Facebook feed; news sites visited
Endline survey subsample	Endline survey, approximately two months after baseline	17,635 (47% of participants who completed the baseline survey)	Political opinions; affective polarization

This table describes the main sample and subsamples analyzed along with the data sources, the number of participants, and the main outcomes. The subsamples and data are described in Section II.C. The outcomes are described in Section II.D.

Table 2: Balance Table, Liberal and Conservative Treatments

Variable	Mean			Difference		
	Sample N=37,494	US	FB Users	Control - Lib.	Control - Cons.	Cons. - Lib.
Baseline Survey						
Ideology (-3, 3)	-0.61	0.17		0.01	0.01	0.00
Democrat	0.38	0.35	0.30	0.01	0.00	0.01
Republican	0.17	0.28	0.21	-0.01	0.00	-0.01
Independent	0.37	0.32	0.35	-0.00	-0.00	-0.00
Vote Support Clinton	0.53			-0.00	-0.00	-0.00
Vote Support Trump	0.26			0.00	-0.00	0.01
Feeling Therm., Rep.	29.07	43.06		0.11	0.25	-0.13
Feeling Therm., Dem.	46.99	48.70		0.40	0.46	-0.06
Difficult Pers., Rep. (1, 5)	3.13			0.02	0.00	0.02
Difficult Pers., Dem. (1, 5)	2.39			-0.00	0.01	-0.01
Facebook Echo Chamber	1.18		1.12	-0.00	-0.00	0.00
Follows News	3.35	2.42		0.01	0.01	-0.00
Most News Social Media	0.18	0.13		-0.00	0.00	-0.00
Device						
Took Survey Mobile	0.67			-0.01*	-0.00	-0.01*
Facebook						
Female	0.52	0.52	0.55	-0.01	-0.00	-0.00
Age	47.69	47.30	42.86	0.22	-0.13	0.35
Total Subscriptions	474			5.15	9.04	-3.89
News Outlets Slant (-1, 1)	-0.18			0.00	0.00	0.00
Access Posts, Pre-Treat.	0.98			0.00	0.01***	-0.00**
Attrition						
Took Followup Survey	0.47			0.03***	0.03***	-0.00
Access Posts, 2 Weeks	0.92			0.00	0.01**	-0.01**
Extension Install, 2 Weeks	0.05			0.00	-0.00	0.00
F-Test				1.20	0.89	1.05
P-Value				[0.21]	[0.64]	[0.39]

This table presents descriptive statistics, along with the difference between participants assigned to each treatment arm. *Vote Support* is the share of participants who voted for or preferred the candidate. *Difficult Pers.* is whether participants find it difficult to see things from Democrats' / Republicans' point of view. *Facebook Echo Chamber* is whether the opinions participants see about government and politics on Facebook are in line with their views always or nearly all the time (3), most of the time (2), some of the time (1), or not too often (0). *Follows News* is whether participants follow government and politics always (4), most of the time (3), about half the time (2), some of the time (1), or never (0). *Total Subscriptions* is the number of Facebook pages participants subscribed to in baseline. *News Outlets Subscriptions* is subscriptions to pages of leading news outlets. *News Outlets Slant* is the slant of news outlets subscriptions. F-tests are calculated by regressing the treatment on the pre-treatment variables, with missing values replaced with a constant and an indicator for a missing value. Data sources for the US and Facebook population are specified in Appendix C.4.1. *p<0.1 **p<0.05 ***p<0.01.

Table 3: Compliance with the Treatments

	(1)	(2)
Cons. Treat., Cons. Ideology	0.513*** (0.008)	
Lib. Treat., Cons. Ideology	0.349*** (0.008)	
Cons. Treat., Lib. Ideology	0.541*** (0.006)	
Lib. Treat., Lib. Ideology	0.623*** (0.006)	
Know Slant		0.230*** (0.006)
Outlet Ideology, Abs. Value (Std. Dev.)		-0.047*** (0.003)
Ideological Distance (Std. Dev.)		-0.083*** (0.002)
Controls	X	X
Observation Unit	Ind.	Ind. * Outlet Offered
Observations	36,728	97,937

This table estimates the association between participants' characteristics and compliance with each treatment arm. In column (1), the dependent variable is whether the participant subscribed to at least one offered outlet and the independent variable is the interaction of participant's ideological leaning and her treatment assignment. The reference group is the control group where there are no compliers. In column (2), the data is pooled such that each observation is a participant and an outlet offered. The dependent variable is whether the participant subscribed to the outlet. The independent variables are based on the outlet's perceived ideology where ideology is measured on a 7-point scale from extremely liberal to extremely conservative with an additional option of 'do not know'. *Ideological Distance* is the standardized difference between the participant's self-reported ideology and the outlet's perceived ideology. Both regressions control for age, age squared, gender, and the set of potential outlets defined for a participant, and column (2) also controls for outlet fixed effects. Column (1) use robust standard errors and column (2) clusters standard errors at the individual level. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Online Appendix - Social Media, News Consumption, and Polarization: Evidence from a Field Experiment

Ro'ee Levy

A Data Collection and Processing

A.1 Surveys

A.1.1 Recruitment Ads

The Facebook ads recruiting participants to the baseline survey mentioned that a research survey was conducted by Yale University and that participants could win Amazon gift cards (Appendix Figure A.13). One version of the ad suggested that the survey was about politics and the other suggested that it was about American society.¹

Most participants were recruited through ads targeting all Facebook users living in the US who are over 18 years old. Using a Facebook Pixel, the ads targeted Facebook users who were more likely to begin the survey. A subset of the ads targeted conservatives or moderate individuals who are often under-represented in Internet samples. Since the majority of participants took the survey on a mobile phone, an additional subset of ads focused on desktop users, to ensure that a large enough sample of participants will be offered an option to install the Chrome extension. A very small minority of users seemed to have a technical issue when taking the survey using the iOS operating system and therefore iOS users were excluded from the target audience once this was discovered (the sample still contains many iOS users). While the survey was open and participants could share the link or ad with anyone, the vast majority of participants probably entered the survey as a result of the ad.²

¹I do not find evidence for heterogeneous effects on political opinions or affective polarization by the type of ad.

²I provided participants with a slightly modified link to the baseline survey after they completed the survey, and asked them to use this link if they wish to share the survey. Only 0.57% of participants entered the survey using this link. Any individual exposed to an ad could also share the ad or the link that appears in the ad with other individuals. Approximately 95% of exposures to the ads during the recruitment period were directly due to a sponsored ad appearing in one's Facebook feed and not due to someone sharing the ad.

A.1.2 Baseline Survey

The baseline survey took place from early February to mid-March 2018. 40,504 responders took the survey and reached the screen where the intervention occurs. Of those, 37,494 are included in the final sample. Responders are excluded from the final sample for the following reasons: missing information on outlets the responder subscribes to either because the responder did not provide permissions to access that data or since the data was not collected properly in real-time (2.38%); the responder already subscribed to too many of the outlets such that it was not possible to define four potential liberal outlets and four potential conservative outlets (4.01%); technical issues with the Qualtrics survey which prevented some data from being collected (0.90%); taking the survey a second time (0.01%); responding carelessly (0.12%). Careless responders are defined as responders who completed all survey questions until the intervention exceptionally quickly (in under three minutes where the median time was eleven minutes) and responders who did not answer at least half of the closed-ended, non-required questions, or who did not answer any question on the final page before the intervention. Finally, to slightly reduce the number of outlets, alternative outlets which are defined as potential outlets for fewer than 20 participants are excluded from the experiment, along with the participants for which these outlets were defined as potential outlets. This removes fewer than 0.1% of participants from the baseline sample.

A.1.3 Endline Survey

Participants were invited to the endline survey between mid-April and early June 2018. Participants were mostly recruited to the survey using emails and Facebook ads.³ To match endline survey responses with baseline survey responses, participants were asked to log in to the endline survey through Facebook or supply an email address. I match endline responses based on the following criteria: email address the survey invitation was sent to, Facebook id, email address entered in the survey, combination of zip code, first and last name if the combination is unique, and combination of first and last name if the combination is unique. 98.73% of endline responses were matched with baseline responders.

17,635 participants are included in the endline survey subsample. If the same individual took the endline survey more than once, uncompleted surveys are excluded. If multiple observations still exist, only the first response is included for the individual. Overall, 0.41% of valid matched responses were excluded as duplicates. 0.02% of responses were also excluded for taking the survey carelessly when the survey was completed exceptionally quickly (spent less than 20 seconds per survey page, compared to a median time of 67 seconds).

³A small share of participants was recruited through an invitation in the browser extension or a Facebook notification.

A.2 Facebook Data on Subscriptions and Posts Shared

I collect data on outlets participants subscribed to (pages “liked”) and posts they shared using a Facebook app, which provides an interface between a Facebook account and the survey.⁴ The data allowed me to customize the survey by ensuring participants are not offered outlets they already subscribed to and including questions about the potential outlets. The app was approved through the standard Facebook review process.

I include in the analysis the following types of shared posts: link, note, status, and video. I focus on these posts since they are more likely to contain political content relevant to the experiment. In some cases, the outlets offered to participants published posts that contain only an image with text (for example, Fox News published posts with quotes related to the news without an accompanying link or video). These posts are defined as photos and are excluded from the analysis. Therefore, the effect I find on the number of posts shared as a result of the experiment is probably slightly lower than the actual effects. When estimating an effect on posts shared, I control for baseline posts shared in the eight weeks before the intervention, when that data exists.

I match posts participants shared with leading outlets based on the Facebook page which published the original post. If a post is not matched with any Facebook page, I determine the slant of the post based on the domain of a link included in the post. For outlets offered in the experiment, I expand the list of domains in the Bakshy, Messing and Adamic (2015) dataset to decrease measurement error. For each outlet, I create a list of relevant domains by checking which domains were shared by the Facebook page associated with the outlet and including the most dominant domains and any other domain related to the outlet. For example, I associate both “huffpost.com” and “huffingtonpost.com” with HuffPost.

If a link refers to a short alias, created by URL-shortening services such as tinyurl.com, it cannot be directly matched to an outlet based on the domain. Therefore, each URL is first converted to the final re-directed URL before being matched to the list of domains.

I also observe participants’ gender and age on Facebook. I define participants’ age as 2018 minus their birth year and replace any age above 90 with missing.

A.3 Extension Data

I collect data on the Facebook feed and browsing behavior using the Google Chrome browser extension. Participants who took the survey on a computer using Chrome were offered to install the extension in the baseline survey before the intervention. In exchange for installing the extension for at least 48 hours, participants could choose between receiving a \$5 gift card, participating in a lottery with a \$200 gift card, or receiving a copy of the study results.

⁴To minimize measurement error, data from the app was collected using several methods, including code running in the background of the baseline survey, a web service, and multiple scripts that ran for the duration of the experiment.

A.3.1 Browsing Behavior

I observe news sites participants visited when the extension was installed. News sites visited are matched to outlets based on their domain. A news site is determined to have been visited through Facebook if the website visited appeared in the participant's Facebook feed in the 20 minutes proceeding to the visit.⁵ I exclude URLs that were visited for less than one second before another URL was visited. If a URL is visited more than once within a 20-minute window, only the first visit is included. When estimating an effect on browsing behavior, I control for baseline browsing behavior in the eight weeks before the intervention, when that data exists.

A.3.2 Facebook Feed

I observe posts appearing in participants' Facebook feeds when participants have the extension installed and use their computer mouse to scroll down the Facebook feed. I do not observe posts unless they appear on the participants' screen. While the extension was designed to work with Google Chrome, it can also work with similar browsers and a very small number of users installed it on alternative browsers, such as Vivaldi.

I assign posts appearing in participants' Facebook feeds to outlets using the following hierarchy:

1. The post was created by a leading news outlet (e.g., a post by the New York Times)
2. The post shared a post created by a leading news outlet (e.g., a friend shared a post by the New York Times).
3. The post includes a link to a leading news outlet (e.g., a friend shared a New York Times link). If the post shares no link, but the text of the post contains a link, I use that link instead. I first convert all links to their final re-directed URL.

I exclude posts where I cannot observe whether the post is shared by a page or a friend (these posts could be sharing content from other Facebook features such as a Facebook Game or Town Hall, they comprise less than 1% of posts in my sample).

In my data, I cannot precisely identify whether a post is sponsored or organic. Instead, I use two techniques to identify ads. First, I assume that any post seen by at least two participants who did not subscribe to the post's page is sponsored. Second, I assume that any post that appeared more than twice in a participant's feed for at least two participants is sponsored. Facebook's algorithm usually does not show the same post many times to the same user, however, advertisers can choose to maximize impressions and thus may show the same post repetitively. When determining whether a post is sponsored, I assume that two posts from the same page with the same text

⁵The time window used is not particularly important. If a 5-minute window is used the number of sites determined to have been visited through Facebook in the two weeks following the intervention decreases by less than 3%, and if a 60-minute window is used, the number of sites increases by less than 3%.

are the same post, even if they have a different id, since advertisers can use two separate posts to run identical advertisements.⁶

While these criteria are far from perfect, they do seem to identify many ads. For example, based on my classification, the top ten words that are most likely to appear in posts identified as ads, compared to organic posts are: “get, now, free, new, today, just, time, one, us, help”. In contrast, the top ten words most likely to appear in organic posts are: “trump, president, one, people, new, school, just, gun, like, now.”⁷

A.4 Leading News Outlets

The list of leading news outlets is based on a dataset of domains constructed by Bakshy, Messing and Adamic (2015). The authors use Facebook’s internal data and classify links as hard or soft news. Hard news articles are related to issues including national news, politics, or world affairs, while soft news includes issues such as sports and entertainment. The alignment of each website is determined according to the self-reported ideology of Facebook users who share hard news links from the website. While many of the sites in the list are traditional news outlets, such as washingtonpost.com, others are more partisan organizations, such as occupydemocrats.com

I exclude from the dataset the following popular websites which are not directly related to news: Amazon, Barack Obama, The White House, Twitter, Vimeo, Wikipedia, and YouTube. I also exclude MSN and AOL since these sites are aggregators of a wide variety of content, they may serve as homepages, and they are often visited for reasons not related to news consumption (Peterson, Shared and Iyengar, 2019). I merge websites that appear twice in the dataset, with and without a web reference, into one entry. For example, washingtonexaminer.com and www.washingtonexaminer.com are merged, with the slant defined as the mean slant of the two entries. After processing the data, the list of leading outlets contains 487 websites.

A.5 Comscore Data

The Comscore Web Behavior Database Panel is a subset of Comscore’s opt-in Media Matrix Panel, which is weighted to represent the US Internet population. Each observation includes a unique machine (computer) id, which I assume represents an individual, although it is possible that multiple individuals use the same machine. When combining data for multiple years, I assign each individual the zip code in the last year for which data exists.

When classifying the referral channel through which a news site was visited, the referring channel is defined as social if the referring domain is one of the following: “facebook.com”, “live.com”, “t.co”, “reddit.com”, “pinterest.com”, “youtube.com”, “linkedin.com”, “twitter.com”,

⁶I make this assumption when the text is at least 20 characters long.

⁷The terms exclude stop words along with the words http, can, said, see.

"tumblr.com", "instagram.com". I classify any referral domain that includes the word google (e.g. "google.com" or "google.co.uk") as a search domain along with the following domains: "yahoo.com", "bing.com", "ask.com", "duckduckgo.com", "searchencrypt.com", "searchlock.com", "searchincognito.com", "search.com", "searchprivacy.co", "safesear.ch", "myprivatesearch.com", "netfind.com". I classify a site as visited directly if there is no referral domain or if the referral domain is the same domain as the domain visited.

B Additional Details on Empirical Strategy

B.1 Segregation Measures

This section describes the isolation and segregation measures in more detail, along with other measures which are presented in Appendix Tables A.7 and A.9.

B.1.1 Isolation

Isolation is the difference between the mean share of conservatives that conservatives are exposed to in the outlets they visit and the mean share of conservatives that liberals are exposed to. Exposure is defined as the share of conservatives browsing the websites among all the site's visitors.

$$Isol = \sum_{i \in \{C_i\}} WeightAmongCons_i * ConsExposure_i - \sum_{i \in \{L_i\}} WeightAmongLib_i * ConsExposure_i$$

where $WeightAmongCons_i$ is the share of outlets visited by individual i among all outlets visited by conservatives, $\{C_i\}$ is the set of conservative individuals, $\{L_i\}$ is the set of liberal individuals, and $ConsExposure_i$ is exposure to conservatives by individual i . Exposure can be calculated as the average share of conservatives among all outlets visited by individual i . To prevent a small sample bias, the average share does not include the visits by i :

$$ConsExposure_i = \sum_j \frac{Visits_{ij}}{Visits_i} * \frac{Cons_j - Visits_{ij}}{Visits_j - Visits_{ij}}$$

where $Visits_{ij}$ is the number of visits of individual i to outlet j and $Visits_i$ is total visits by individual i , so $\frac{Visits_{ij}}{Visits_i}$ is the weight of outlet j for individual i . $Visits_j$ is total visits to site j and $Cons_j$ is total conservative visits to site j , so $\frac{Cons_j - Visits_{ij}}{Visits_j - Visits_{ij}}$ is the share of conservatives visiting outlet j excluding individual i .

B.1.2 Segregation

Segregation is defined as the scaled standard deviation of partisan news exposure. This can be interpreted as the expected square distance between the slant of news sites visited by two random participants in the sample (Flaxman, Sharad and Rao, 2016):

$$Seg = \sqrt{2} * std.dev(Slant_i)$$

where $Slant_i$ is the mean slant of outlets visited by individual i . The slant of outlet j is based on Bakshy, Messing and Adamic (2015) and normalized to the unit interval (by adding one and dividing by two).

B.1.3 Absolute Value of Slant

To measure the extremity of news consumption, I calculate the absolute value of mean consumption slant as:

$$AbsSlant = \sum_i \frac{|Slant_i|}{N}$$

where $Slant_i$ is the mean slant of outlets visited by individual i and N is the number of individuals in the sample. The slant of outlet j is based on Bakshy, Messing and Adamic (2015) such that a middle-of-the-road outlet has a slant of zero, a completely conservative outlet has a slant of 1 and a completely liberal outlet has a slant of -1.

B.1.4 Congruence

I define congruence as exposure to more extreme content matching the consumer's ideology:

$$Congruence = \sum_i \frac{(Slant_i * IdeoLeaning_i)}{N}$$

where $Slant_i$ has the same definition as in the previous measure and $IdeoLeaning$ is defined as 1 for a conservative participant and -1 for a liberal participant. N is the number of individuals in the sample for which an ideological leaning can be defined.

B.1.5 Share of Counter-Attitudinal News

To determine the share of counter-attitudinal news, I divide news sites into five quintiles: very liberal, liberal, moderate, conservative, and very conservative (Bakshy, Messing and Adamic, 2015). I define pro-attitudinal news as conservative and very conservative news consumed by a conservative, or liberal and very liberal news consumed by a liberal. Counter-attitudinal news is

conservative and very conservative news consumed by a liberal, or liberal and very liberal news consumed by a conservative. Finally, the share of counter-attitudinal news is defined as the share of counter-attitudinal news among all pro- and counter-attitudinal news.

$$ShareCounter = \sum_i \frac{\frac{\sum_j (IdeoLeaning_i = SlantGroupOutlet_j)}{\sum_j SlantGroupOutlet_j \in \{-1, 1\}}}{N}$$

where $SlantGroupOutlet_j = 1$ if outlet j is conservative or very conservative and $SlantGroupOutlet_j = -1$ if outlet j is liberal or very liberal.

B.2 Pre-Analysis Plan

The main outcome and hypotheses tested in this study were pre-registered in the AEA RCT Registry.⁸ The analysis deviates from the pre-analysis plan in two important ways. First, I use equal weights when constructing the indices, while the plan states that the weights will be determined by the inverse of the covariance between the outcome measures (Anderson, 2008). This method is not used since it generates negative weights. With negative weights, the interpretation of an index is less clear. For example, the question on President Trump’s approval rating received a negative weight which means that *ceteris paribus*, a participant who has a more favorable opinion on Trump would be considered more liberal.

Appendix Table A.18a estimates the effect on the political opinions index using equal weights in column (1) and inverse-covariance weights in column (2). This method does not cleanly generate weights for individuals with missing outcomes. In column (3), weights from column (2) are renormalized to sum to one for participants with missing outcomes, an index is then created for each participant by weighting the standardized outcomes, and finally, the index is standardized with respect to the control group. Since the inverse-covariance method generates negative weights, columns (4) and (5) repeat the analysis with negative weights replaced with zero and the weights renormalized accordingly. While there is some variation in the results, the most straight-forward comparison is between columns (1) and (5). These columns focus on the same participants and do not include negative weights. In column (5), the effect of the conservative treatment is slightly larger but still small in magnitude and not statistically significant.

Appendix Table A.18b shows that the effect on affective polarization is robust to using inverse-covariance weights.

The second important deviation from the pre-analysis plan is that the polarization index originally included five attitudinal measures and three behavioral measures, while only the attitudinal measures are analyzed in this paper. The behavioral measures were based on a question in the endline survey asking participants whether they would “like” or share a post stating that “In

⁸AEA RCT Registry Trial 0002713.

seeking truth, you have to get both sides of a story.” The primary behavioral outcome is composed of an index of the following measures: did participants state they will share the post, did participants state they would “like” the post, did participants actually share the post. However, it was not possible to analyze the actual behavior of a large share of participants partly due to the unexpected Cambridge Analytica scandal, which led many individuals to revoke access to the posts they share. Furthermore, the behavioral measure turned out not to measure polarization well. While a measure of polarization should typically be correlated with partisanship, there was almost no correlation between being partisan and the behavioral outcomes.⁹

Column (1) of Appendix Table A.19 shows that the effect is still significant when using all eight variables in the polarization index.¹⁰ Column (3) measures the effect only on the behavioral outcomes (for most participants data not exist on whether posts were shared so this index is mostly based on the self-reported survey answers). The effect of the treatments is small and not statistically significant. While this result does not change the conclusions regarding affective polarization, it is interesting to note that exposure to counter-attitudinal outlets does not affect participants’ self-reported willingness to share or like a post on seeking both sides of a story.

When processing and analyzing the data, I made various other minor changes compared to the pre-analysis plan, include the following. In the plan, I stated that I will estimate the results excluding the first two days after the intervention. Instead, I estimate the results for each week or month separately. The plan states that the regression will control for the randomization block and for whether the participant used the iOS operating system. I exclude the iOS variable for simplicity (this does not affect the primary endline survey results). I do not control for the randomization blocks (strata) since due to attrition, some blocks have only one or two respondents instead of the original three respondents. When controlling for the block, I am only able to analyze a subset of participants. The results for that subset are essentially the same with and without controlling for strata. I do not report raw or adjusted p-values for each index component of the political opinions and affective polarization measures, as I do not focus on the individual components. Instead, I present each component visually in appendix figures.

In the pre-analysis plan, ideological leaning is defined first by self-reported ideology and then by party affiliation. I prefer using party affiliation as the main variable defining ideological leaning to make the study comparable to other papers, which tend to focus on party affiliation (Druckman and Levendusky, 2019). The results are robust to the original definition. I also control for ideological leaning in the primary endline survey regressions. In contrast to the plan, I do not present several demographic variables in the balance table since they suffer from post-treatment

⁹The correlation between the behavioral polarization measures and the absolute value of a baseline scale of partisan affiliation (where 0 is no party identification, 1 is leaning toward a party, 2 is identifying with a party and, 3 is strongly identifying with a party) is only 0.04-0.06. The correlation between the affective polarization measures and partisan affiliation is 0.22-0.46.

¹⁰The effect when all eight variables are used to construct a polarization index is smaller in index points than the effect when the five attitudinal measures are used. When standardizing the indices with respect to the control group, the effects are similar since the index created when using all eight variables has less variation in the control group.

bias and do not impute them since I already have rich survey and social media data. Finally, the pre-analysis plan states that a political knowledge index will be created. Since I do not focus on political knowledge, I instead analyze separately the effect on each political knowledge primary outcome in Appendix C.6. While the results are easier to interpret when analyzed separately, an index would not change the qualitative conclusions of the section.

B.3 Controls

To increase power, when estimating the effect on political opinion and affective polarization, I control for a set of pre-registered covariates. I control for self-reported ideology, party affiliation, approval of President Trump, ideological leaning, age, age squared, gender. Age and gender are included in the Facebook data provided when participants log in to the survey and the remaining covariates are based on the baseline survey. Self-reported ideology is a nominal variable with seven ideological options from very liberal to very conservative and an option for participants who have not thought much about this. Party affiliation is a nominal variable with seven affiliation options ranging from strong Democrat to strong Republican along with an option of “other party”. Approval of Trump is a nominal variable with four options ranging from strongly disapprove to strongly approve.

When estimating the effect on political opinions, I also control for the following baseline survey questions: feeling toward President Trump (0-100 integer); worry about illegal immigration (nominal variable with the options not at all, only a little, fair amount, great deal); does the participant believes Mueller is conducting a fair investigation (nominal variable with the options yes, no, do not know), and whether the participant thinks Trump has attempted to obstruct the investigation into Russian interference in the election (nominal variable with the options yes, no, do not know).

When estimating the effect on affective polarization, I also control for the baseline values of the *feeling thermometer* and *difficult perspective* measures (defined in Section II.D.2).

In all regressions, if a covariate includes missing values, the missing values are coded to a constant and an additional dummy control is added to the regression indicating whether a value is missing. Regressions testing for heterogeneous effects also control for each participant’s potential outlets since individuals who were assigned the alternative outlet may have different characteristics than individuals who were assigned the primary outlets.

C Additional Analysis

C.1 Survey Purpose

At the end of the baseline survey, participants were presented with the following question: "If you had to guess, what would you say is the primary purpose of this study?" Appendix Table A.20

shows the most common phrases participants mentioned according to their treatment assignment. Unsurprisingly, participants understood that the study is on media and politics, as most questions focused on these topics and the consent form stated that this is the topic of the study. Among the most common phrases, there are not many substantial differences between the treatments.

Appendix Table A.21 presents the phrases with the largest differential usage between the treatment arms and the control group. While participants in both the pro- and counter-attitudinal treatments mentioned terms such as “echo chamber” and “social media” more often than the control group, probably due to the text of the intervention encouraging participants to “Like” Facebook pages, the differences between the two treatment arms in the usage of these terms is small. When comparing the pro- and counter-attitudinal treatments to each other, almost no substantial differences stand out. One exception is that a small share of participants in the counter-attitudinal treatment thought the purpose of the survey was to get them to like liberal Facebook pages. These participants probably were not pleased with the experimenter trying to “push liberal” content (that was not the actual purpose of the experiment, of course) and therefore it is unlikely that they expressed opinions aligned with these outlets to make an impression on the experimenter. In any case, while these phrases represent a relatively large difference between the treatments, they are not mentioned often.

Overall, this section suggests that participants in the counter-attitudinal treatment did not perceive the experimenter’s expectations substantially differently than participants in the pro-attitudinal treatment. This conclusion does not rule out that experimenter effects played a role in some of the results. It is possible, for example, that participants in the pro- and counter-attitudinal treatments understood that the study attempts to analyze the effect of news outlets on political opinions, they remembered which outlets they were offered, and tried in the endline survey to convey attitudes more similar to the outlets offered (e.g., a more positive opinion toward the Republican Party if they were offered conservative outlets). However, at least it is unlikely that differential expectations of the experimenter’s objective are driving the main results.

C.2 Analysis of the Content Participants Engaged With

In this section, I show that the most common content participants engaged with as a result of the intervention is political. I analyze the posts from the subscribed outlets that participants were exposed to in their feed, links in the posts that they visited, and posts they shared using three methods. First, I show the most common phrases mentioned in the posts. Second, I define certain terms as political and analyze the share of political posts. Third, I analyze the section and outlet where each article appeared based on the URLs appearing in the posts.

An important challenge in this analysis is that the posts affected by the treatment cannot be cleanly identified. For example, participants in the control group visited the news sites of their potential counter-attitudinal outlets approximately 1.70 times in the two weeks following the intervention,

while participants in the counter-attitudinal treatment visited these websites approximately 1.34 additional times (as shown in Figure 6). While the participants were affected by the treatment, I cannot identify which of their visits to counter-attitudinal news sites would have occurred in a counterfactual with no intervention. I focus on posts affected directly by the intervention by analyzing only posts shared by pages participants subscribed to in the experiment (excluding suspected ads). While this decreases the likelihood of including posts that participants would have engaged with without the intervention, it does not cover the entire effect of the intervention. For example, participants often visited the websites of the offered outlets indirectly, even when they did not observe the specific link to an article in their feed (as shown in Figure A.5).

Throughout this section, I focus on the eight weeks following the intervention to increase the number of data points. To reduce variability in the text analyzed, I include in the analysis only posts from the eight primary outlets and first two alternative outlet that were offered to participants. This excludes less than 3% of posts participants were exposed to.

Before discussing the results, an important caveat is in order. This section is descriptive and its purpose is to show what content participants engaged with according to whether the outlets they were offered were pro- or counter-attitudinal. When comparing the content shared by liberals who subscribed to liberal outlets (pro-attitudinal) with content shared by conservatives who subscribed to liberal outlets (counter-attitudinal), I am *not* estimating the causal effect of the treatments, as the compositions of the two groups compared are different by definition.

C.2.1 Most Common Phrases

Appendix Table A.22 shows the most common phrases mentioned in posts participants were exposed to in their feed, in posts with links participants visited, and in posts shared by participants. I first remove punctuation, terms that appear in only one outlet, media-related terms or terms that were likely to be covered mostly by specific outlets (e.g., “write” or “New York”), and then stem the words appearing in the posts.¹¹

The most common phrases participants were exposed to are political and are usually related either to President Trump, the aftermath of the Parkland school shooting, or the Mueller investigation. The phrases appearing in posts participants clicked are similar to the phrases in posts participants were exposed to.

The posts shared should not be directly compared to the posts participants were exposed to or clicked since the data is based on two different subsamples. Regardless, it is clear that posts shared are often political even when participants shared posts in the counter-attitudinal treatment. However, the response to scandals may be heterogeneous. For example, liberals are more likely to share

¹¹In addition to stop words, I remove the following terms: bit, breaking news, can, comment, fox friend, fox news, http, https, journal, last week, new york, new york time, news, nyt, opinion, said, say, times, wall street journal, washington post, write, write the editori board, wsj, year old.

articles mentioning Robert Mueller in both the pro- and counter-attitudinal treatments. Similarly, liberals in the liberal treatment are more likely to share articles mentioning Stormy Daniels and conservatives in the conservative treatment are more likely to share articles mentioning Hillary Clinton.

C.2.2 Share of Posts Mentioning Political Words

Focusing on the most common words allows us to understand which topics were most prominent but does not provide a complete analysis of the posts, especially if there is a lot of variability in the posts' content. In this subsection, I use a simple measure to determine a lower bound for the share of political posts. I define a post as political if it contains terms related to political figures ("biden, bolton, carson, clinton, devos, kushner, manafort, mccabe, mcconnell, michael cohen, obama, pelosi, pence, pruit, tillerson, trump"), political parties ("conservative, democrat, dnc, gop, liberal, republican, the left, the right"), political institutions ("congress, elect, politic, senate, vote, white house") or political issues ("ar 15, daca, gun control, gun law, gun right, immigration, mass shooting, nra, parkland, sanctuary city, sanctuary state, school shooting, tax cut, walkout"). I search for the terms in the post's text, its URL, and any commentary on the post if it is shared.¹²

Remarkably, more than half of the posts observed, clicked, and shared, are political. This is probably a lower bound for the actual number of political terms since posts including the terms I mentioned are almost always political but there are other political posts not captured by these terms (e.g., posts about race relations, gender issues, climate change and additional posts about gun legislation that do not include a unique term that can be clearly identified as political).

Appendix Figure A.14 shows that participants in the pro-attitudinal treatment were generally more likely to engage with political posts. However, the difference between the pro- and counter-attitudinal treatments is surprisingly small with one notable exception. Among liberals who shared posts from liberal outlets they were offered, 68% of posts were political, compared to 41% among conservatives who shared posts from the offered liberal outlets.

Still, it may be surprising that a large portion of the counter-attitudinal posts shared by participants was political. Why do participants share these posts? Anecdotally, there seem to be various reasons. Some posts are written by moderate columnists in a counter-attitudinal outlet (e.g., William A. Galston at the Wall Street Journal), others focus on rare bipartisan topics (e.g., a bill against sex trafficking), or report topical news without expressing strong opinions. In other cases, the posts may tackle issues where the outlet does not completely share the party's line, or where the participants may not agree with the party (e.g., conservatives who oppose the NRA's positions). There were also cases where participants share the posts with a negative comment, even though these are less common than might be expected. Finally, in a few cases, participants admitted they are sharing posts from outlets they usually would not share. This suggests that typically

¹²Specifically, for shared posts I search for political terms in the message, description, and link fields.

participants did not start sharing partisan news completely supporting the other side, but they may have shared articles with more nuanced positions in counter-attitudinal outlets.

C.2.3 Outlets and Sections

Instead of determining the posts' topics based on words in the post, I can analyze the content participants engaged with using the outlets' own classification of their articles. Most outlets classify articles into sections, such as News, Business, and Arts, and mention the sections on their website, the website's HTML, or the URL. I determine the section associated with a post based on analyzing the website associated with the post. This method is not perfect. MSNBC usually does not classify articles and videos into sections and Slate often creates short links for its URLs which were no longer available when I determined the link's section. Still, the advantage of this method is that it relies on internal decisions by the outlets, who should know their content best.

Appendix Figure A.15 shows the most common outlets and sections participants were exposed to. The figure mostly reflects the different preferences of participants when subscribing to outlets. Liberals mostly avoided "liking" Fox News when it was offered and preferred the Wall Street Journal. They were more likely to already subscribe to one of the primary liberal outlets in baseline, and therefore, more likely to be offered to subscribe to Washington Post, the first alternative liberal outlet.

Appendix Figure A.16 suggests that participants clicked a larger share of posts about culture or arts compared to the share observed in the feed. For example, entertainment articles from HuffPost and cultural articles from the Washington Times are more prominent in this figure. Interestingly, this holds both for participants in the pro- and counter-attitudinal treatments. However, posts with links to politics and national news are still most likely to be clicked in both treatments.

The differences between the posts shared by participants are more stark. For example, Appendix Figure A.17 shows that conservatives shared HuffPost articles in the parenting, women, or queer voices sections, while among posts shared by liberals, these sections form a very small minority.¹³ Still, within each outlet, the dominant sections among posts shared are typically the political or national news sections, even in the counter-attitudinal treatment.

C.3 Heterogeneous Effects

In the pre-analysis plan, I stated that I will test for heterogeneous effects based on whether participants are ideological, whether they are in an echo chamber, the openness of participants, and whether they are sophisticated.

¹³Interestingly, almost no articles shared were in the sports section (less than 1% of articles for which a section could be identified).

I define participants as *Ideological* if the absolute value of their self-reported ideology on the 7 point scale (from -3 for very liberal to +3 for very conservative) is above or equals the median.

I use two measures of being in an echo chamber. The variable *Echo Chamber* is whether the answer to "Thinking about the opinions you see people post about government and politics on Facebook, how often are they in line with your own views" is above or equals the median. *Seen Counter Att.* is whether the share of potential counter-attitudinal outlets, among all potential outlets, participants reported seeing in their feed in baseline is above or equals the median.

I measure whether a participant has an *Open Personality* according to whether her average agreement with the following statements is above or equals the median: "I see myself as open to new experiences, complex" and the reverse values of "I see myself as conventional, uncreative." The questions are based on Gosling, Rentfrow and Swann (2003). I define participants as *Certain* in their opinions if their answer to "Generally speaking, how certain are you of your political opinions?" is above or equals the median.

I define participants as *Sophisticated* if they answered one of the following questions correctly: "Suppose 110 members of a local government voted on an infrastructure bill. The bill passed by a margin of 100 votes. How many members voted against the bill", "Suppose the number of US citizens on the internet doubles every month. If it took 48 months for the entire US population to have internet access, how many months did it take for half the population to have internet access". These questions are based on the Cognitive Reflection Test (Shane, 2005).

In addition to the pre-registered tests, I explore the effect of several additional moderators. *Most News Social Media* is whether participants reported getting most of their news about government and politics through social networking sites. Participants have *High News Subscriptions* if their baseline number of subscriptions to pages of news outlets on Facebook is above or equals the median. Participants are considered *Exposed to Outlets* if their self-reported exposure to posts from the eight potential outlets in baseline is above or equals the median. Participants are considered to be *Familiar with Slant* if the distance between their perceived slant of the potential outlets and the average perceived slant by participants with the same self-reported ideology is below the median. Participants are considered to *Follow the News* if their answer to "how often do you pay attention to what's going on in government and politics?" is above the median. Participants are considered to have a *High Feeling Thermometer Difference* if the difference between their feeling toward their own party and the opposing party is above or equals the median. Finally, participants are considered *Conservative* if their ideological leaning is conservative, *Older* if their age is above or equal to the median age, and *Female* if they identify in Facebook as female.

When analyzing heterogeneity in the effects of the pro- and counter-attitudinal treatments, I do not distinguish between heterogeneity due to differences in the participants' ideology and heterogeneity due to differences in the outlets offered. For example, if conservatives are affected more by the pro-attitudinal treatment, that could be due to conservatives being more persuadable or because Fox News is more persuasive than New York Times.

Appendix Figures A.18 and A.19 estimate heterogeneous effects on subscribing to outlets, exposure to posts from outlets, and visiting the outlets' websites. Each row represents a separate regression estimating the effect of interacting the pro- or counter-attitudinal treatment with the specified variable, where the reference group is the control group.¹⁴ A higher value means individuals were more likely to engage with the pro- or counter-attitudinal potential outlets as a result of the pro- or counter-attitudinal treatment, respectively.

Ideological individuals were more likely to subscribe to pro-attitudinal outlets and less likely to subscribe to counter-attitudinal outlets. Participants who were more certain in their opinions, and who follow the news were also less likely to subscribe to counter-attitudinal outlets. Similarly, ideological participants, along with participants following the news and participants who were more polarized in baseline, were less likely to visit these outlets. Finally, participants who subscribed to more outlets in baseline were more likely to subscribe to counter-attitudinal outlets. Interestingly, even though they subscribed at higher rates, they were *less* likely to be exposed to these outlets in their feed as a result of the intervention, probably since there is more competition for space in their feed.

Appendix Figure A.20 estimates heterogeneous effects on the primary endline survey outcomes and shows that the effect on political opinions is mostly homogeneous (i.e., most participants were not persuaded by the treatments). The right panel of Appendix Figure A.20 does not show strong heterogeneous effects on affective polarization according to most covariates tested. The strongest heterogeneous effect found is based on the baseline feeling thermometer measure for affective polarization. The effect on affective polarization is weaker among participants who were more polarized in baseline. However, this result is significant at the 10% level and the results are not adjusted for multiple hypothesis testing, and therefore more research is needed to explore heterogeneity in affective polarization.

C.4 Reweighting for National Representativeness

C.4.1 Data Sources

To reweight the sample to match the US population, I use the following data sources. The medium where Americans get most of their news is based on the Pew American Trends Panel Wave 23 (November-December 2016). All other US data is based on the 2016 American National Election Survey (ANES). The estimates are based on pre-election ANES questions, besides vote or support for a presidential candidate, which is based on the post-election survey.

In Table 2, I also present demographics for Facebook users. Data on whether the opinions Facebook users see about government and politics on Facebook are in line with their views is based on

¹⁴The results of most heterogeneous effects are similar when estimating all the heterogeneous effects on either political opinions or affective polarization simultaneously in one regression.

a question in the Pew American Trends Panel Wave 1 (March-April 2014) asked among respondents who pay attention to posts about government and politics on Facebook. All other data on Facebook users is based on the 2018 Pew Core Trends Survey.

C.4.2 Analysis

In this section, I reweight the sample to match the national population using the entropy weighting procedure (Hainmueller, 2012). I match the following subset of control covariates: self-reported ideology (mean value on a scale of 1-7), the share of participants identifying as Democrats, Republicans, and Independents, the difference between the participants feeling toward their party and the opposing party, age, and the share of females. For the feeling thermometer, self-reported ideology, age, and gender covariates, missing variables are first replaced with the mean value (less than 5% of observations are missing for each of these variables). When analyzing the effects of the pro- and counter-attitudinal treatments, I compare the sample to the US population for which an ideological leaning can be defined and use those means to reweight the sample.¹⁵

Appendix Tables A.23 and A.24 show that reweighting the sample does not change the main conclusions of the study. The effect on the slant of posts participants were exposed to increases slightly. The effect on sites visited, posts shared, political opinion, and affective polarization remain essentially the same, although the confidence intervals are wider. These tables should be interpreted with caution. It is likely that even after reweighting, the sample is still different than the national population on unobservables or covariates not used when reweighting the sample. Still, the tables show that it is unlikely that an effect on affective polarization is only found because the survey sample is more liberal or more polarized than the rest of the population.

C.5 Predicted Treatment Effect for the Full Baseline Sample

The previous section reweighted participants to match the US population. In this section, I predict the main treatment effect for the entire baseline sample. While the baseline sample is not nationally representative, such an estimation provides several advantages. First, it estimates the same results among a larger group of participants that are more representative than the extension and endline survey subsamples, using a large set of Facebook and survey covariates. Second, it alleviates concerns that differential attrition by some observable characteristics is driving the results.

I first estimate heterogeneous effects on the slant of posts observed, the slant of news sites visited, the political opinions index, and the affective polarization index. The effects on media engagement are estimated in the extension subsample and the effects on self-reported opinions and attitudes

¹⁵I include respondents who identify or lean toward one of the parties, who define themselves as liberal or conservative, or who voted, intended to vote or preferred Donald Trump or Hillary Clinton, according to the ANES pre-election survey. Overall, 94% of respondents in the ANES survey are included.

are measured in the endline survey subsample.¹⁶ I exclude the control group in these estimates so the interpretation is the effect of the conservative treatment on conservative media consumption and conservative opinions, compared to the liberal treatment, or the effect of the pro-attitudinal treatment on polarization, compared to the counter-attitudinal treatment. I estimate heterogeneous effects using causal forests (Wager and Athey, 2018). The intuition behind causal forests is that one part of the sample is used to determine how to split each tree and another part is used to estimate heterogeneity. If the same sample was used for both processes, heterogeneity would be overestimated due to overfitting, as the sample would be split according to the covariates that happen to predict heterogeneous effects in this particular sample.

I use a large set of covariates including almost all close-ended baseline survey questions and data from Facebook on the age, age squared, and gender of the participant, the number of pages liked by the participant in baseline, and the number of pages the participant liked in 2017. In addition, I include covariates for whether each of the outlets in the experiment could have been potentially offered to the participants and whether the participant liked a set of popular pages on Facebook (for example, one variable is whether the participant liked The Beatles on Facebook). I include all pages liked by at least 10% of participants in baseline. In total, 255 covariates are used. I then use these covariates to predict the ITT effect among all participants in the baseline sample.

Appendix Table A.25 shows that the results predicted among the entire baseline sample are very similar to the results found among the subsamples of participants who completed the endline survey or installed the Chrome extension for at least two weeks. Based on the analysis of heterogeneity throughout this paper, the fact that the effects on opinions and attitudes are stable is not surprising, as the effects on the primary outcomes are generally homogeneous and the differences between participants in the baseline and endline surveys are not dramatic.

While these results are reassuring, two caveats should be noted. First, I control for many observable variables, but there could be unobservables differentiating the subsamples. Second, when estimating heterogeneous effect in the extension subsample, I cannot control for one important difference between the groups - the device with which the survey was taken - since participants could only install the extension when taking the survey on a computer using Google Chrome.

C.6 Effects on Knowledge

While this paper focuses on persuasion and polarization, the endline survey includes several questions related to political knowledge. The two primary measures of political knowledge are self-reported familiarity, measured according to whether participants reported hearing of news events and political figures, and accurate political knowledge, measured according to participants' answers to several questions on recent events. For some questions, participants were expected to

¹⁶I do not analyze the effect on posts shared because the access posts subsample already includes a large share of the baseline sample.

gain knowledge when assigned to the liberal treatment (heard of Michael Cohen, heard about the Stephon Clark shooting, believed the Russian government tried to influence the 2016 elections, believed a wall is not being built at the US-Mexico border) and for other measures, the conservative treatment was expected to have an effect (heard of Louis Farrakhan, heard about a controversial speech by Hillary Clinton in India, believed Trump is not a criminal target of the Mueller investigation, believed Trump's tax cuts would increase most people's income).

Appendix Table A.26 shows that the treatments had little to no effect on the knowledge outcomes. The coefficients of interest are the effects of the liberal treatment on liberal outcomes and conservative treatment on conservative outcomes. Most of the coefficients are small in magnitude and not statistically significant.

Appendix Table A.27 tests whether there is no substantial effect on knowledge because the treatment did not affect exposure to the topics the endline survey focused on. The table uses the extension data to estimate the effect of the treatments on posts appearing in the participants' social media and shows that the intervention affected all four self-reported familiarity outcomes (Michael Cohen, Stephon Clark, Louis Farrakhan, and the Hillary Clinton speech).¹⁷

The results presented in this section suggest that while the slant of one's social media feed can determine the news events an individual is exposed to on social media, that exposure does not necessarily affect their political awareness of topics. One possible explanation is that individuals consume news also outside their social media feed. In any case, this result should not be interpreted as definitive evidence of a null effect. Participants were asked questions about specific issues and answers to knowledge questions could be driven by motivated reasoning.

C.7 Exposure to Posts From the Offered Pro- and Counter-Attitudinal Outlets

In this section, I provide more details on the decomposition exercise in Section VI, analyze several alternative decompositions, and test whether there is a gap in exposure to pro- and counter-attitudinal posts within outlets.

C.7.1 Decomposition Calculations

I include in this analysis participants in the pro- and counter-attitudinal treatments for which I can observe posts in the Facebook feed in the two weeks following the intervention and for whom at least one post is observed. Overall, the sample includes 521 participants in the pro-attitudinal treatment and 538 participants in the counter-attitudinal treatment.

¹⁷Posts are defined as referring to Michael Cohen, Louis Farrakhan, or the shooting of Stephon Clark if they include the terms "michael cohen", "louis farrakhan" and "stephon clark," respectively. Posts refer to Hillary Clinton's speech in India suggesting that many white women voted for Trump since they took their voting cues from their husbands if they include the words "clinton," "vote," and either "india" or "husband."

I define the number of posts from counter-attitudinal outlets observed in the counter-attitudinal treatment as:

$$S_C * A_C * U_C$$

where S_C is the mean number of new subscriptions to the offered counter-attitudinal outlets; A_C is the effect of the algorithm determining the share of posts in the feed from the subscribed counter-attitudinal outlets among all the posts in the feed (formally defined later in this section); and U_C is the total number of posts observed in the feed in the counter-attitudinal treatment. I define the number of posts observed in the pro-attitudinal treatment as:

$$S_P * A_P * U_P = (S_C + S_\Delta) * (A_C + A_\Delta) * (U_C + U_\Delta)$$

I then decompose the difference in exposure to four separate expressions as described in Equation 3. To estimate S_Δ and U_Δ , I use the following regressions:

$$TotalSub_i = S_\Delta ProTreat_i + \varepsilon_i$$

$$TotalPosts_i = U_\Delta ProTreat_i + X_i + \xi_i$$

where $TotalSub_i$ and $TotalPosts_i$ are the number of offered outlets the participant subscribed to and the total number of posts observed, respectively. These regressions are presented in columns (1) and (2) of Appendix Table A.28. X_i controls for Facebook usage before the intervention to increase precision.

To estimate the effect of subscribing to a post on exposure, I pool the two groups of potential outlets such that for each participant there are two observations: one observation with the four potential pro-attitudinal outlets and one observation with the four potential counter-attitudinal outlets. I calculate the share of posts the participants observed from each group of outlets among the total number of posts from all sources the participant observed in the two weeks following the intervention. I only include posts shared directly by the outlet to isolate any effect of friends sharing specific posts. I use the share of posts as the outcome variable instead of the total number of posts since users may observe more posts from pro-attitudinal outlets due to increased Facebook usage, and I account for that effect separately. A_C and A_Δ are estimated using the following regression:

$$SharePosts_{ij} = A_C * Sub_{ij} + A_\Delta * Sub_{ij} \times Pro_{ij} + \delta * Pro_{ij} + v_{ij} \quad (1)$$

where $SharePosts_{ij}$ is the share of posts participant i observed from group j , Sub_{ij} is the number of outlets participant i subscribed to from group j . Pro_{ij} is whether the outlets in the group are pro-attitudinal. I instrument for Sub_{ij} and $Sub_{ij} \times Pro_{ij}$ with $Offer_{ij}$ and $Offer_{ij} \times Pro_{ij}$, where $Offer_{ij}$ is whether participant i was offered outlets from group j in the intervention. This regression is presented in column (3) of Appendix Table A.28. Conceptually, it can be easier to think of this regression as two separate regressions. One regression includes only the potential counter-attitudinal outlets, and measure the effect of subscribing to an outlet on exposure to the outlet

(A_C). I exploit the fact that for some participants the counter-attitudinal outlets were offered and for others they were not offered. In a second regression, I repeat this exercise for the potential pro-attitudinal outlets. A_Δ is the difference between the coefficients.

C.7.2 Alternative Decompositions

Appendix Figure A.21 presents the decomposition exercise using several alternative estimations. The x-axis is the gap in exposure to posts from the pro- and counter-attitudinal outlets, in the two weeks following the intervention. Most of these specifications lead to similar results, although I am often underpowered to detect precise effects. The first row of the figure is the primary specification shown in Figure 10. The second row adds fixed effects for the potential outlets defined for each participant. This assures that the estimates are derived from comparing participants who could have been offered the same set of outlets. The rest of the decompositions are described below.

Exclude Unsubscriptions Participants in the counter-attitudinal treatment may observe fewer posts due to their decision to unsubscribe from the offered outlets. Since they initially subscribed to the outlet, this could be accounted for as an algorithmic effect. In the third row of Appendix Figure A.21, only subscriptions lasting at least two weeks are defined as subscriptions (this estimation only includes participants for which I observe two weeks of subscription data). The results do not change substantially.

Exclude Suspected Ads In the primary decomposition, I assume that Facebook’s algorithm determines whether participants observe posts from outlets they subscribe to. This typically holds for organic posts. However, participants also observe sponsored posts (ads) which are different in several important aspects. First, they can appear in a user’s feed even if she did not subscribe to the outlet. Second, the placement of sponsored posts can be determined by the advertiser. For example, an outlet can promote posts to a subset of users who subscribed to its Facebook page. This means that part of the effect attributed to the algorithm may result from the behavior of advertisers.¹⁸ When excluding suspected ads, the gap between exposure to pro- and counter-attitudinal outlets slightly decreases. This suggests that ads target users whose ideology matches the outlet they subscribe to. Still, even when ads are excluded, the gap between the two groups of outlets remains large and the decomposition does not change substantially.

¹⁸Even with sponsored posts, the algorithm may still play an important role. For example, advertisers can target a broad array of users and pay for each click on a post. This creates an incentive for Facebook to place the posts among users who are likely to click them, and thus the incentives in determining where to place sponsored posts can be similar to the incentives when placing organic posts.

Reweight Based on Compliance The effect of the algorithm is estimated using two IV estimators, and thus its causal interpretation relies on the assumption that there is no essential heterogeneity (Heckman, Urzua and Vytlačil, 2006). Otherwise, the difference between exposure to posts, conditional on subscriptions, in the pro- and counter-attitudinal treatments might be due to the combination of heterogeneity in the effect of subscribing to outlets and selection into compliance, and not due to differing effects of subscribing to pro- and counter-attitudinal outlets. In the fifth row panel of Appendix Figure A.21, I re-weight the IV estimators, such that participants predicted to comply receive a lower weight. I first calculate the probabilities of compliance with the pro- and counter-attitudinal treatments, by regression compliance on the following covariates using a logit regression: age, female, self-reported ideology, party (dummy variables for Democrat, Republican, and Independent), and the difference between the participant's feelings toward her party and the opposing party. I then predict the probability of compliance for each participant and define the participant's weight as the inverse of the predicted probability.

The figure shows that reweighting the compliers does not change the result substantially. The reweighted estimates measure the treatment effect under the conditional effect ignorability assumption (Angrist and Fernandez-Val, 2013; Aronow and Carnegie, 2013). This assumes that conditional on the covariates (the compliance score), subscribing to outlets has the same average treatment effect for compliers on non-compliers. There could still be essential heterogeneity based on other variables differentiating the compliers, but at least this suggests that the result does not stem from differences in compliers and heterogeneous effects by ideology or baseline affective polarization, for example. The result is similar to the main estimate not because the effect is homogeneous, but rather because the compliers are not dramatically different from non-compliers in both treatments.

Reweight to Match Population Demographics In the sixth row of the figure, I reweight the participants to match population means on the same set of variables mentioned in the previous section using the entropy weighting procedure. Reweighting decreases the gap in the number of posts observed. When analyzing the results separately for conservatives and liberals, I find that the algorithm's tendency to increase exposure to matching news outlets is driven by the liberals in my sample (I am underpowered to estimate this result precisely) and that could explain the decreased gap in exposure when reweighting the results.¹⁹ Still, there remains a substantial gap in exposure to pro- and counter-attitudinal posts even after reweighting the participants.

Excluding Facebook Usage The effect on Facebook usage is only marginally significant. In the seventh row of Appendix Figure A.21, I assume that the exposure gap only stems from subscriptions and the platform algorithm, and exclude the usage dimension. For this decomposition, I change the calculation of A in equation 1, and instead of estimating the effect on the share of posts

¹⁹The difference between liberal and conservatives could be due to the ideology of participants or differences in the outlets offered.

in the feed, I estimate the effect on the number of posts observed by participant i from outlets in group j .

Decomposition Over Time In the final two rows of Appendix Figure A.21, I decompose the gap in exposure for the first and second week after the intervention. I use the same estimate for subscriptions in both weeks but calculate exposure to posts and Facebook usage according to each week's specific activity. The overall gap in the number of posts is greater in the first week, but this reflects the fact that participants were generally exposed to more posts from the offered outlets in the first week. The relative difference between pro- and counter-attitudinal posts is greater in the second week (approximately 140% more pro-attitudinal posts) compared to the first week (106%). The effect associated with subscriptions becomes smaller over time and the effect associated with the algorithm slightly increases. This suggests that Facebook's algorithm learns from participants' behavior that they prefer pro-attitudinal content. However, the effect of the algorithm is still strong in the first week suggesting that either the algorithm learns very quickly (e.g., based on engagement with the first posts from an offered outlet shown to a participant) or that the algorithm uses other baseline information (such as subscriptions to other outlets) to determine that participants are more interested in pro-attitudinal content.

C.7.3 Differential Exposure to Articles Within an Outlet

To estimate whether participants were exposed to news more likely to match their opinions within an outlet, I focus on the subset of articles that were shared on Facebook or Twitter by at least one member of Congress in January-November 2018. I define the slant of an article according to the mean first dimension of the DW-Nominate score of congress members who shared the article (Jeffrey et al., 2020).²⁰

I find that in general conservative participants are exposed to more conservative articles on Facebook, even when controlling for the outlet. This is not surprising as a conservative is likely to have more conservative friends, who are likely to share more conservative articles within an outlet. However, when I focus only on posts shared by the eight potential outlets defined for each participant, I do not find any correlation between the slant of the posts and consumers' ideologies. This suggests that Facebook's algorithm does not lead to conservatives being supplied with more conservative articles, *within* the set of posts shared by an outlet. It also suggests that conservatives and liberals were exposed to similar content from the outlets they subscribed to in the intervention, conditional on posts from the outlet appearing in their feed.

²⁰The list of the Facebook pages of congress members is based on the Congress Members Project (<https://github.com/unitedstates/congress-legislators>). Based on this list, I collected all posts shared by congress members in 2018. The list of tweets shared by congress members is from the Tweets of Congress Project (<https://github.com/alexlitel/congresstweets>). The datasets were downloaded in December 2018.

D Interpretation

How should we interpret the fact that the intervention affected attitudes toward parties, while political opinions remained stable? In this section, I compare two frameworks explaining affective polarization and examine which is most consistent with the data.

Consider the following model: consumer i 's prior on state k of the world is $\theta_{ik} \sim (\theta_{ik}^0, \frac{1}{h_{ik}})$, where θ_{ik}^0 is the consumer's initial belief and h_{ik} is the precision of the belief (the consumer's certainty). I extend classic media persuasion models by introducing the concept of affective polarization and assuming that a consumer's political opinion, γ_i , is a weighted average of K beliefs:

$$\gamma_i = \sum_{k \in \{1..K\}} w_{ik} \theta_{ik} \quad (2)$$

where $w_{ik} \in \{0, 1\}$ is the weight consumer i places on belief k when determining her political opinion. A weight can be thought of as the priority the consumer places on a specific belief. For example, a consumer's support for a climate bill can depend on two beliefs: the consumer's belief on whether the bill will decrease or increase emissions and the belief on whether the bill will increase or decrease electricity prices. A liberal may place a positive weight only on the effect on emissions and a conservative may place a positive weight only on the effect on prices.²¹ A political party uses the same framework and its opinion is a weighted average of various beliefs.

Outlet j receives signal s_{jk} on the state of the world: $s_{jk} \sim N(\theta_k^*, \frac{1}{h_{jk}})$, where θ_k^* is the true state of the world and h_{jk} is the precision of the signal received. Media outlets act as delegates for their consumers by covering issues according to the weights their consumers place on them.²² Therefore, pro-attitudinal outlets cover issues more when $w_{own} > w_{opposing}$ and counter-attitudinal outlets cover issues more when $w_{opposing} > w_{own}$, where w_{own} are the weights used by the individual's own party and $w_{opposing}$ are the weights used by the opposing party. Indeed, Figure 2 suggests that there is substantial differentiation in the topics news outlets cover. Returning to the climate change example, data from the outlets offered in the experiment also demonstrates this differential coverage: for every post from a conservative outlet mentioning the words "environment" or "climate," 1.28 posts mentioned the word "economy," while for liberal outlets, the ratio was 0.43.²³

²¹In the Pew Research Center Political Survey from January 2019, 74% of Democrats stated that the environment should be a top priority for President Trump and Congress in 2019, compared to only 31% of Republicans. On the other hand, 79% of Republicans said the economy should be a top priority, compared to 64% of Democrats (the sample includes respondents leaning toward the Democratic and Republican parties). As a clarifying example for the framework, I intentionally focus on a broad issue, support for climate change policy. Some of the questions forming the political opinions index focus on more specific topics, but the same logic holds. For example, the favorability of the March for Our Lives Movement could depend on participants' belief on whether banning certain weapons will decrease gun violence and their belief on whether the movement will prevent most gun owners from purchasing their preferred guns.

²²Delegation has long been suggested as an explanation for why consumers prefer like-minded news (Suen, 2004; Chan and Suen, 2008).

²³This calculation is based on the ratio between the number of times the words "economy", "climate" and "environment" appeared in the messages of all posts shared by the eight primary outlets and first two alternative outlets between February 15, 2018, and December 31, 2018. Duplicate posts with the same message are excluded.

I assume that consumers exposed to a new outlet update their beliefs in the direction of the outlet. This type of movement is expected if media outlets are biased in their reporting and consumers are naive and do not completely take the bias into account (DellaVigna and Kaplan, 2007).²⁴

A straightforward way to model affective polarization is to define attitudes as a linear function of the distance between the political opinion of party p and a benchmark for the “correct” opinion according to individual i :

$$A_{ip} = g(\gamma_p - \hat{\gamma}_{ip}) \quad (3)$$

where A_{ip} is the attitude of individual i toward party p , γ_p is the political opinion of party p and $\hat{\gamma}_{ip} = \phi(\theta_{i1}, \dots, \theta_{ik}, w_{i1}, \dots, w_{ik}, \theta_{p1}, \dots, \theta_{pk}, w_{p1}, \dots, w_{pk})$, is the benchmark opinion that individual i thinks party p should hold. I consider two benchmark opinions: either individuals use their own opinion as the benchmark or they determine the benchmark opinion based on their beliefs weighted by the weights party p places on the beliefs.

Affective polarization due to political distance: $A_{ip} = g(\gamma_p - \sum_k w_{ik}\theta_{ik})$

Consumers may determine their attitudes toward a party based solely on the distance between their opinion and the party’s opinion, i.e., they use their own opinion as the benchmark for the opinion the party should hold. Without loss of generality, I will focus on the position of a liberal consumer toward the Republican Party ($\gamma_i < \gamma_p$). When the individual’s political opinion changes from γ_i^0 to γ_i^1 due to a change in her beliefs, the following change is expected in her attitude toward party p :

$$\Delta A_{ip} = g(\gamma_p - \gamma_i^1) - g(\gamma_p - \gamma_i^0) = g(\sum_k w_{ik}(\theta_{ik}^0 - \theta_{ik}^1)) \quad (4)$$

According to this theory, increased affective polarization can be explained by ideological divergence (Rogowski and Sutherland, 2016). An update in the consumer’s beliefs should only affect attitudes toward a party through its effect on the consumer’s political opinions. Returning to the climate bill example, a consumer would determine her attitude toward a political party based on the distance between her support for the climate bill and the party’s support for the bill. If a liberal’s support for a bill increases she will develop more negative attitudes toward a party opposing the bill. This theory is not consistent with the experiment since attitudes changed without a corresponding change in political opinions.

²⁴An alternative explanation for why consumers’ posteriors move toward the opposing party when exposed to counter-attitudinal news is that individuals’ priors tend to support their political opinion. In other words, liberals tend to have more liberal priors than the true state of the world and conservatives tend to have more conservative priors. When exposed to counter-attitudinal outlets, liberals and conservatives receive more signals on issues for which they have weak prior and their beliefs move toward the true state of the world.

Affective polarization due to unreasonable opinions: $A_{ip} = g(\gamma_p - \sum_k w_{pk}\theta_{ik})$

Alternatively, the attitude of a consumer toward a party may depend on whether the political opinion of a party is reasonable according to the party's weights. Hence, the benchmark opinion is the opinion the party would hold according to the consumer's beliefs regarding the state of the world, weighted by the weights party p places on those beliefs. In other words, affective polarization increases when consumers cannot rationalize the parties' political opinions and perceive that the party is not adhering to its values.²⁵ The change in affective polarization following an update to the consumer's beliefs is:

$$\Delta A_i = g(\gamma_p - \sum_k w_{pk}\theta_{ik}^1) - g(\gamma_p - \sum_k w_{pk}\theta_{ik}^0) = g(\sum_k w_{pk}(\theta_{ik}^0 - \theta_{ik}^1)) \quad (5)$$

If the consumer and the party place the same weight on beliefs ($w_{pk} = w_{ik}$), there is no difference between the two theories. However, with heterogeneous weights, political opinions and affective polarization may be differentially affected. In the climate bill example, a liberal who believes the climate bill will mitigate emissions and *decrease* consumer prices will support the bill. The consumer will have a negative attitude toward a party opposing the bill since even if the party places a zero weight on decreasing emissions, it should still support the bill. If the liberal is exposed to conservative outlets and learns that the bill is likely to increase prices, she may still support the bill since she places a positive weight only on mitigating emissions but will develop a less negative attitude toward a party that places a positive weight on consumer prices and thus opposes the bill.²⁶

This theory is consistent with the results of the experiment if the consumers updated beliefs on which they place zero weights, but at least one of the parties places positive weights.²⁷ This would result in consumers' political opinions remaining constant, but attitudes toward parties changing.²⁸

²⁵Another way to interpret affective polarization according to this framework is that the consumer attributes malicious motives to the party. Since the consumer infers that the party should have a different political opinion according to its weights and the correct beliefs, she concludes that there is an additional unethical consideration determining the party's stance. For example, the consumer might assume that the party supports a policy because it is corrupt or because the policy will have negative implications for the party's opponents.

²⁶Stone (2020) shows that affective polarization could increase due to limited strategic thinking or a false consensus bias. In the context of this experiment and theoretical framework, a false-consensus bias is similar to consumers having the wrong priors regarding the weights the opposing party places on beliefs. Exposure to counter-attitudinal news allows consumers to learn those weights and thus rationalize the opinions of the opposing party. I focus on beliefs regarding issues and not beliefs regarding the opposing party's weights because I suspect that weights are more likely to be common knowledge. However, both theories are consistent with the results of my experiment.

²⁷It is plausible that as a result of the experiment consumers updated beliefs on which they place zeros weights since they are less likely to have been exposed to counter-attitudinal outlets covering these beliefs. Thus, they are expected to have weaker priors regarding those beliefs. Indeed, Appendix Figure A.4 shows that participants assigned to the counter-attitudinal treatment were more likely to say that they modified their views in the past two months because of something they saw on social media, compared to participants assigned to the pro-attitudinal treatment.

²⁸The stability of political opinions relies on a strong assumption that consumers place zero weights on some beliefs or that they determine their political opinions based on lexicographic orderings of beliefs. This assumption is plausible in certain cases. For example, individuals who do not believe climate change is happening may place a zero weight on

To further test these theories, I analyze the effect of the experiment on participants' attitudes toward the opposing party. If affective polarization is simply a function of political distance, attitudes toward parties will be affected when consumer i updates beliefs on which she places positive weights (Equation 4). Therefore, attitudes toward both parties are more likely to be affected by pro-attitudinal outlets that cover these beliefs. On the contrary, if affective polarization is a function of unreasonable opinions, attitudes toward party p will be affected more by beliefs on which p places positive weights (Equation 5). As a result, pro-attitudinal outlets are more likely to affect attitudes toward one's own party, while counter-attitudinal outlets are more likely to affect attitudes toward the opposing party. Appendix Table A.17 shows that attitudes toward the opposing party are indeed more likely to be affected by exposure to counter-attitudinal outlets, consistent with the theory that affective polarization is due to opinions that are perceived to be unreasonable.

To conclude, there is still limited evidence on whether exposure to pro- and counter-attitudinal news has an effect on affective polarization, let alone an understanding of the channels explaining this effect. I present a parsimonious theory that is consistent with the results: consumers determine their attitudes toward a party based on the distance between the party's opinions and the opinion the party should hold according to the consumers' beliefs and the party's weights. While I provide evidence supporting the theory, there could be other explanations for the change in affective polarization, and more research is needed to pinpoint the precise mechanisms explaining how affective polarization evolves.

E Additional Figures and Tables

whether a climate bill decreases greenhouse gas emissions. More importantly, the logic behind the theory still holds if consumers place a positive but small weight on beliefs. In that case, we would expect political opinions to be slightly affected when those beliefs change, but the effect could still be much smaller than any change in affective polarization.

Figure A.1: Example for the Conservative Treatment Intervention

Following a news or media page is a great way to learn about the news and hear other perspectives. Recently, researchers have suggested that subscribing to random sources can help burst the social media echo chamber.

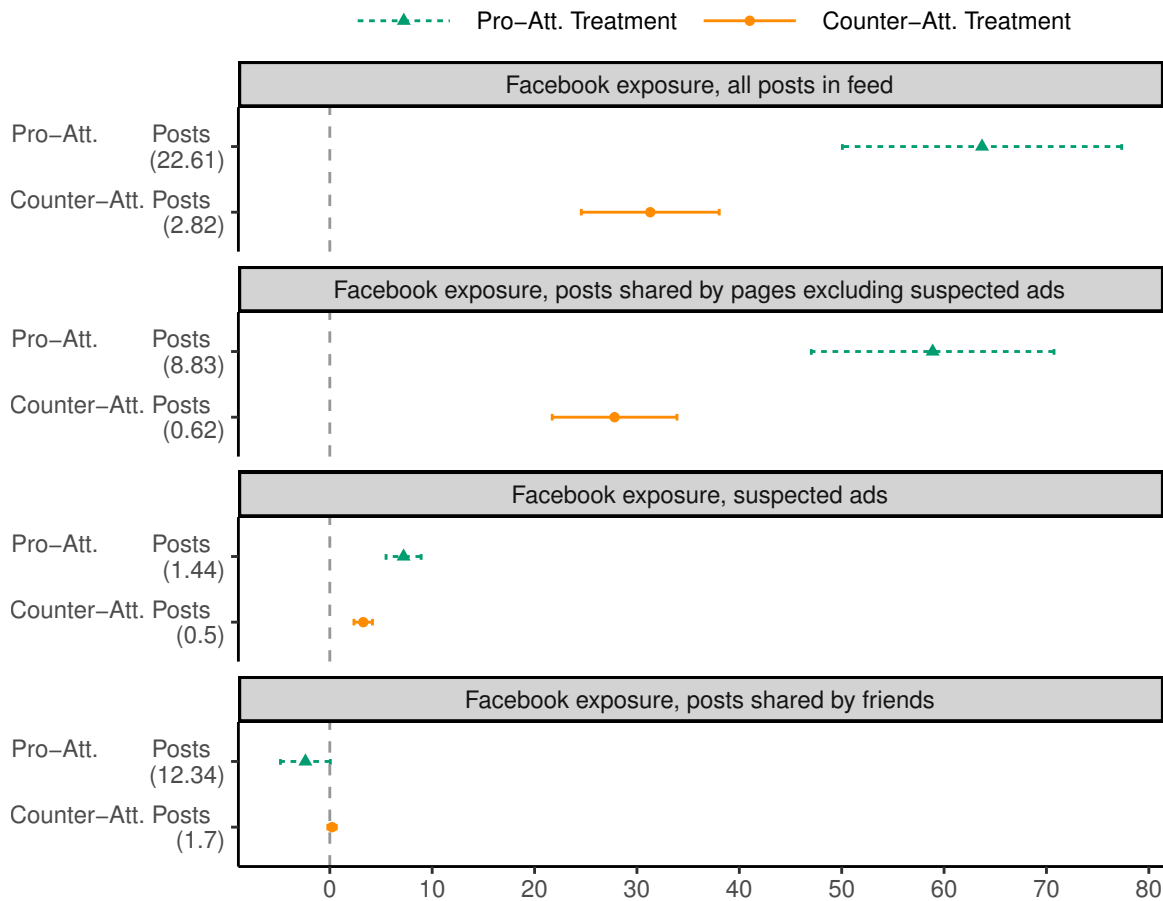
By clicking like below, posts from randomly chosen popular Facebook pages may start appearing in your news feed. **To expand your horizons, please click "Like Page" on 1-4 of the pages below** (Facebook may ask you to confirm the like, you can always unlike the page later).

The pages were chosen randomly and therefore may all represent views you agree or disagree with. In any case, they present an opportunity to diversify your news feed.



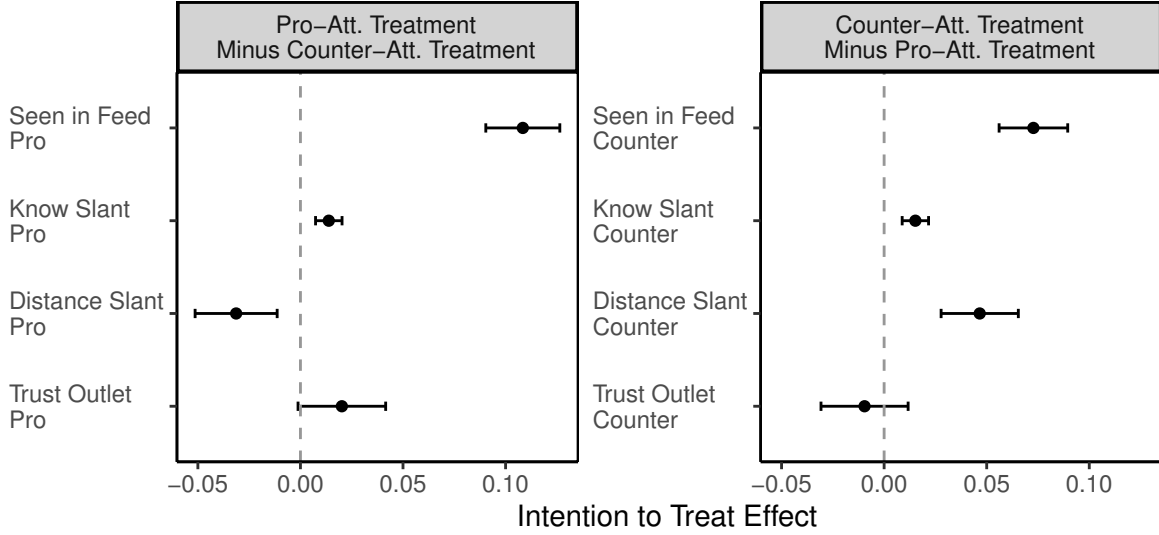
This figure shows a survey page asking participants to subscribe to four conservative outlets. Participants randomly assigned to the conservative treatment, who have not already subscribed to the four primary outlets, were presented with an intervention similar to this figure. The “Like Page” buttons were generated using Facebook’s Page Plugin. The image in the background of each button was automatically updated according to the outlet’s Facebook page, and the order of the outlets was determined randomly.

Figure A.2: Effect of the Pro- and Counter-Attitudinal Treatments on Exposure to the Potential News Sites, by Type of Post



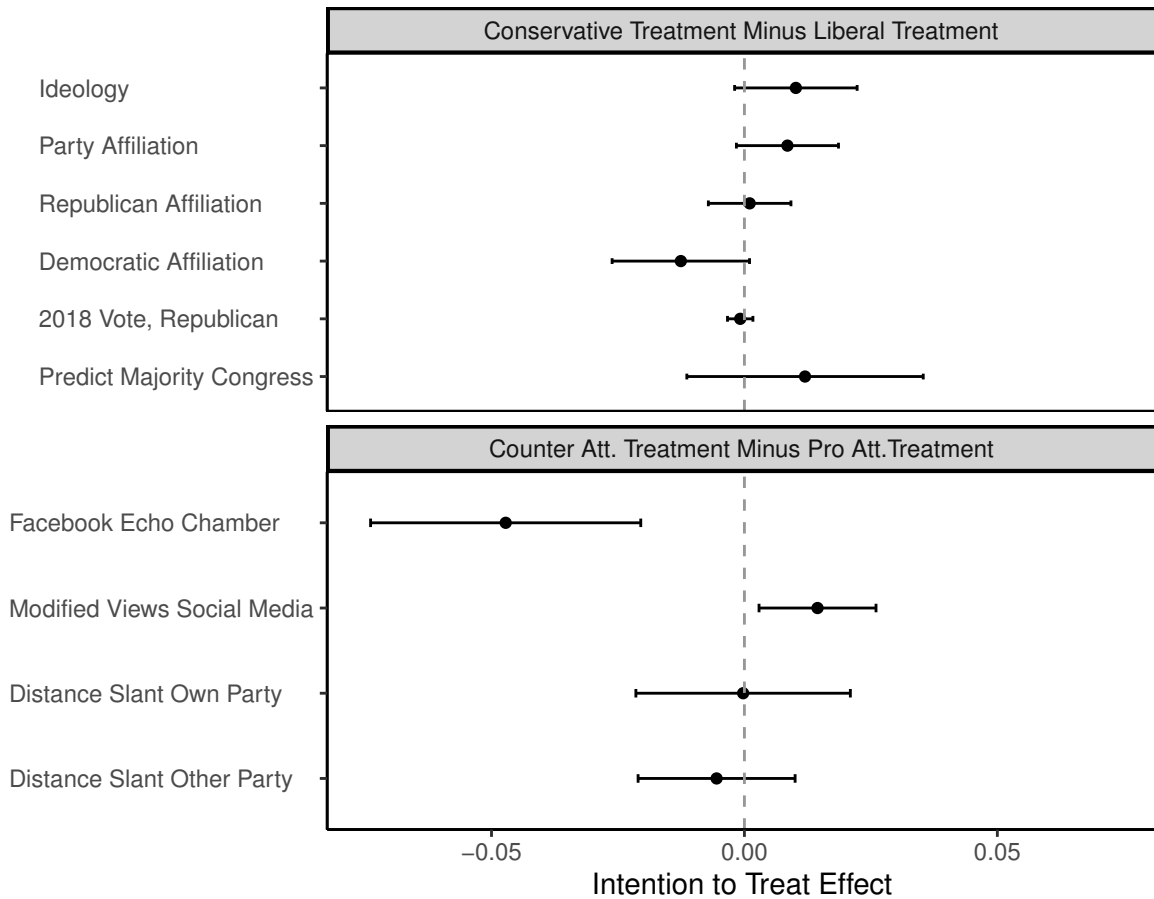
This figure shows the effect of the pro-attitudinal and counter-attitudinal treatments on exposure to posts from the potential outlets in the two weeks following the intervention. The control group mean for each outcome is in parenthesis. The first panel showing total exposure is identical to the second panel of Figure 6. The second panel shows the effect on posts shared by Facebook pages organically. This includes all posts shared by the potential outlets, or other Facebook pages referring to the potential outlets, besides posts which are likely to be sponsored (ads). The third panel shows the effect on exposure to suspected ads related to the outlets. The fourth panel shows the effect on posts shared by Facebook friends. Appendix A.3 explains how ads were identified. Error bars reflect 90 percent confidence intervals.

Figure A.3: Effects on Survey Responses Related to the Potential Outlets



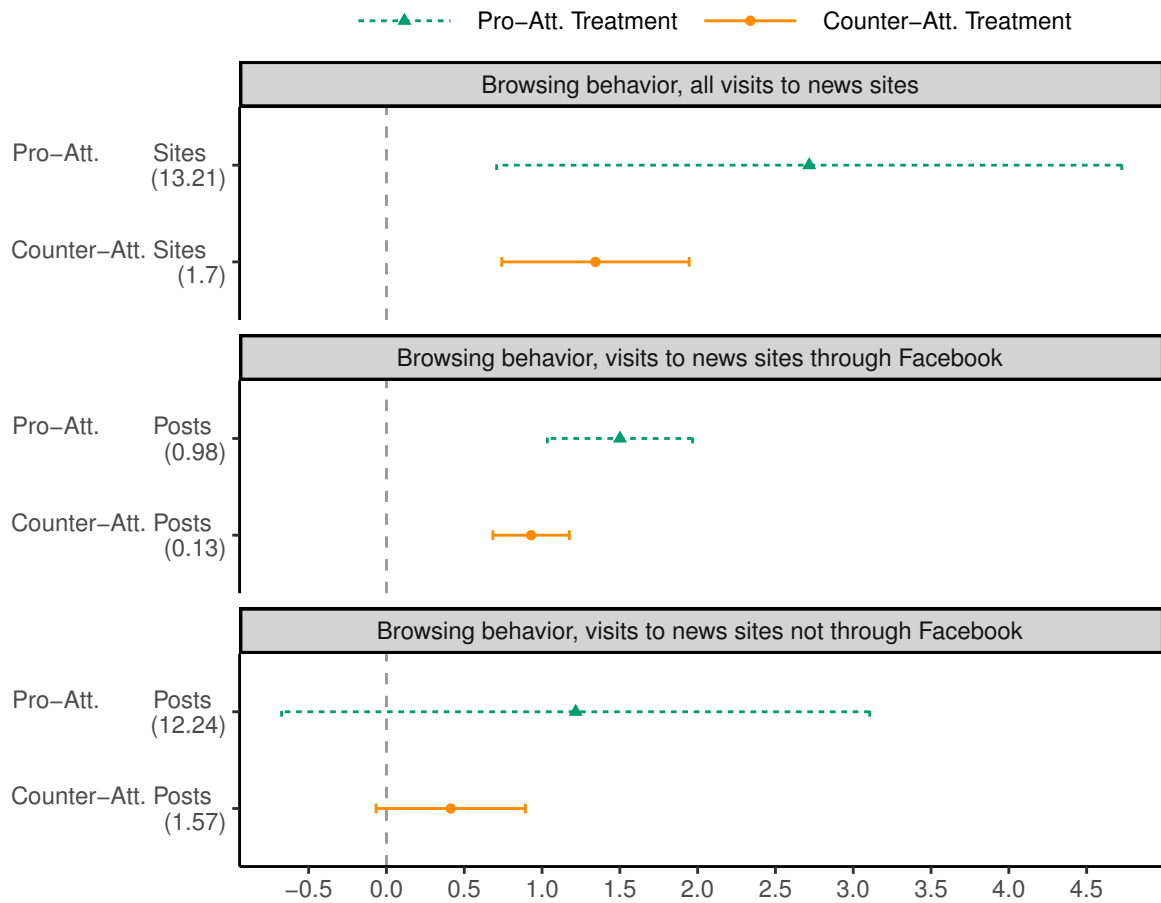
This figure shows the effect of the experiment on attitudes toward the potential outlets. Each row represents a regression pooling the opinions of participants in the endline survey on the eight potential outlets defined for each participant. *Seen in Feed* is whether the participant reported seeing news from the outlets in their Facebook feed over the past week more than five times (3), 3-5 times (2), 1-2 times (1), or reported seeing no posts (0). *Know Slant* is whether the participants did not mark “do not know” when asked what is the outlet’s slant. *Distance Slant* is the difference between the participant’s baseline ideology and the perceived ideology of the outlet. *Trust Outlet* is whether the participant perceived the outlet as very trustworthy (2), trustworthy (1), not trustworthy nor untrustworthy (0), untrustworthy (-1), or very untrustworthy (-2). Non-binary outcomes are standardized by subtracting the control group mean and dividing by the control group standard deviation. The left panel shows the effects of the pro-attitudinal treatment on the pro-attitudinal outlets (the counter-attitudinal treatment is the reference group). The right panel shows the effects on the counter-attitudinal treatment on counter-attitudinal outlets. In addition to the standard controls (Section II.E), the regressions control for baseline outcomes when they exist, outlet fixed effects, and the set of potential outlets defined for each participant. Standard errors clustered at the individual level. Error bars reflect 90 percent confidence intervals.

Figure A.4: Effect of the Treatments on Additional Survey Outcomes



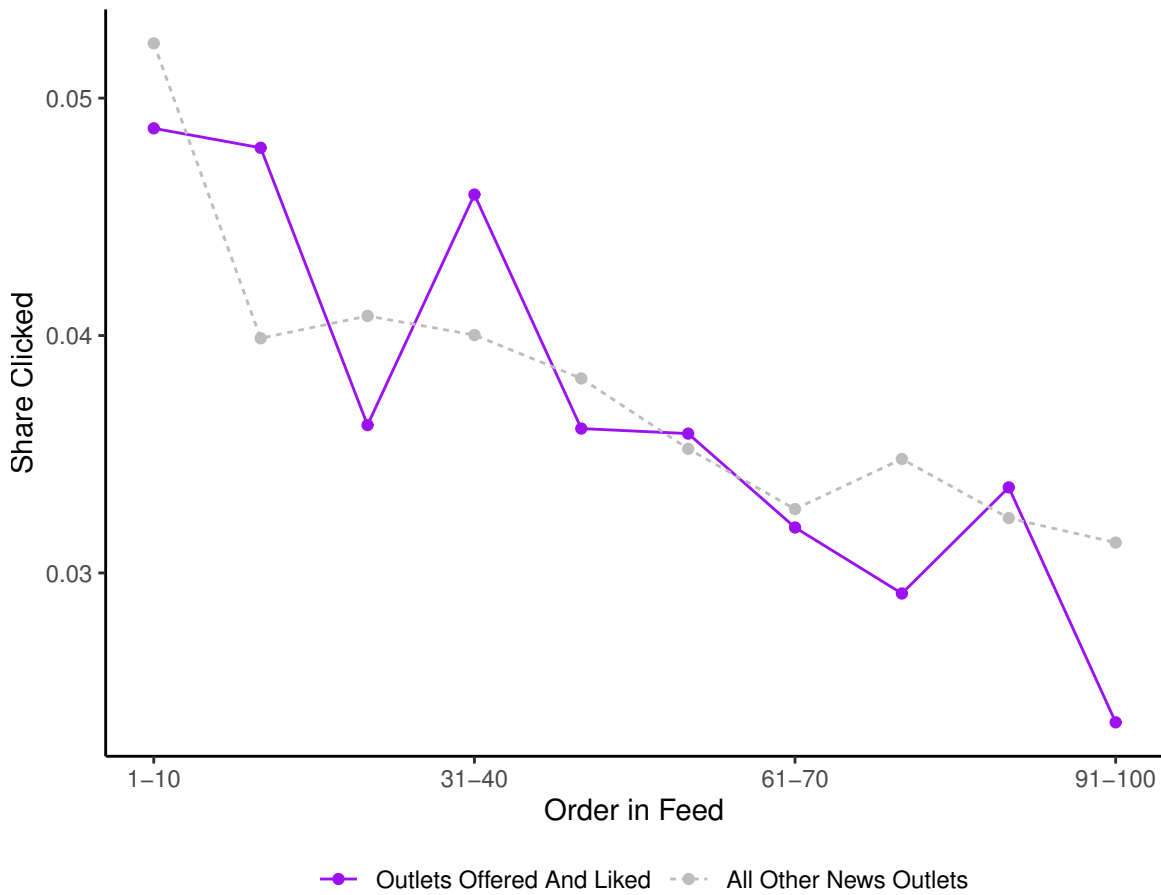
This figure shows the effect of the experiment on additional endline survey outcomes. *Ideology* is self-reported on a 7-point scale. *Party Affiliation* is the party the participant identifies with on a 7-point scale from strong conservative to strong liberal. *Republican/Democrat Affiliation* is whether the participant is a strong Republican/Democrat (3), is a Republican/Democrat (2), leans toward the Republican/Democratic Party (1), or does not identify with both parties (0). *2018 Vote* is whether the participant intends to vote for the Republican Party candidate (1) or the Democratic Party candidate (0) in her district if the election was held the day the survey was taken. *Predict Majority Congress* is the party the participant's predicts will hold the majority of seats in Congress after the 2018 vote: Republican Party (1) not sure (0), or the Democratic Party (-1). *Facebook Echo Chamber* is whether opinions seen about government and politics on Facebook are in line with participants' views always or nearly all the time (3), most of the time (2), some of the time (1), not too often (0). *Modified Views Social Media* is whether the participant modified her views in the past two months about a political or social issue because of something she saw on social media. *Distance Slant* is the difference between the participant's baseline ideology and the perceived ideology of a party. Non-binary outcomes are standardized by subtracting the control group mean and dividing by its standard deviation. In addition to the standard controls (Section II.E), the regressions control for baseline outcomes when they exist. Error bars reflect 90 percent confidence intervals.

Figure A.5: Effects of the Pro- and Counter-Attitudinal Treatments on News Sites Visited, by Source



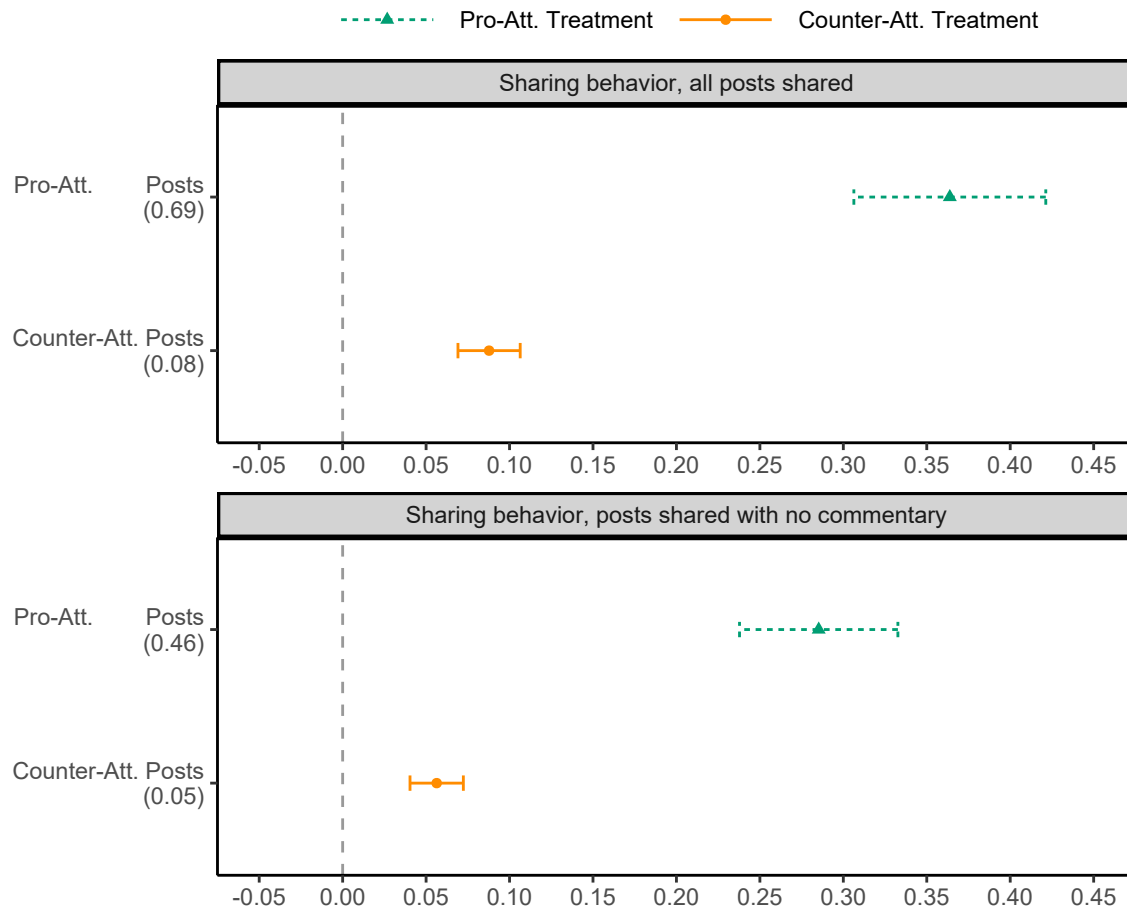
This figure shows the effect of the pro- and counter-attitudinal treatments on total visits to the potential outlets' websites in the two weeks following the intervention. The control group mean for each outcome is in parenthesis. The first panel showing total visits is identical to the third panel of Figure 6. The second panel shows the effect on visits to websites that could be matched with a URL appearing in a Facebook post. The third panel shows the effect on all other visits. Appendix Section A.2 explains how posts were matched with visits to news sites. Error bars reflect 90 percent confidence intervals.

Figure A.6: Share of Links Visited by Order in Feed



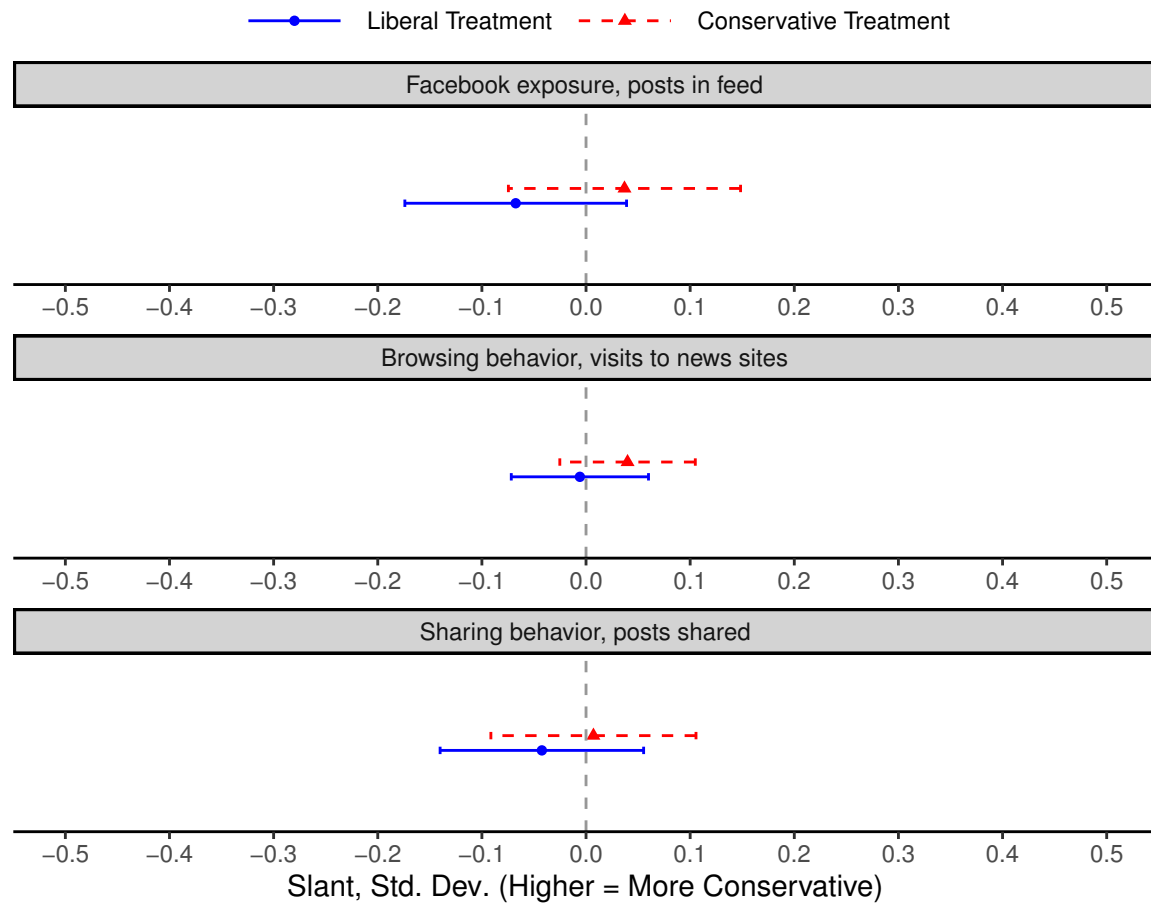
This figure shows the share of links which were visited by participants. The data includes all posts with links from the pages of leading news outlets, excluding suspected ads, in the two weeks following the intervention. To determine the order of posts, a Facebook feed session is defined to begin when a participant views a post on Facebook at least 30 minutes after viewing a previous post. To smooth the results, posts are grouped into groups of ten based on their order. Appendix A.2 explains how posts were matched with visits to news sites and Appendix A.3 explains how suspected ads were identified.

Figure A.7: Effects of the Pro- and Counter-Attitudinal Treatments on Number of Posts Shared, Access Posts Subsample



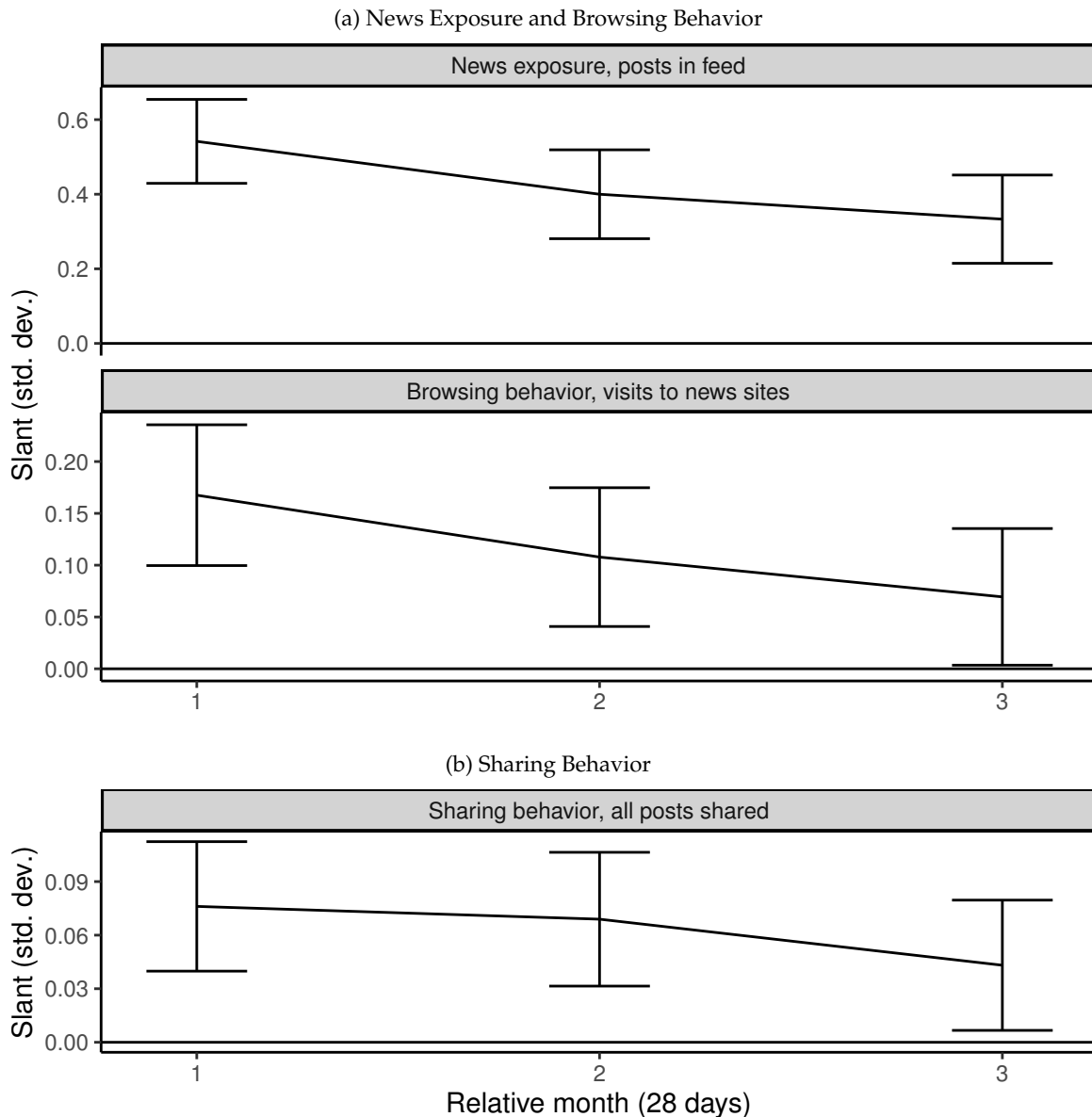
This figure shows the effect of the pro- and counter-attitudinal treatments on the number of posts participants shared from the four potential pro-attitudinal outlets and four potential counter-attitudinal outlets in the two weeks following the intervention. The control group mean for each outcome is in parenthesis. The first panel includes all posts and the second panel includes only posts that were shared without any commentary by the participant. The regressions control for the outcome measure in baseline. The data is from the access posts subsample: 33,532 participants with a liberal or conservative ideological leaning who provided access to their posts for at least two weeks following the intervention. Error bars reflect 90 percent confidence intervals.

Figure A.8: Effect of the Liberal and Conservative Treatments on Slant, Excluding each Participant's Eight Potential Experimental Outlets



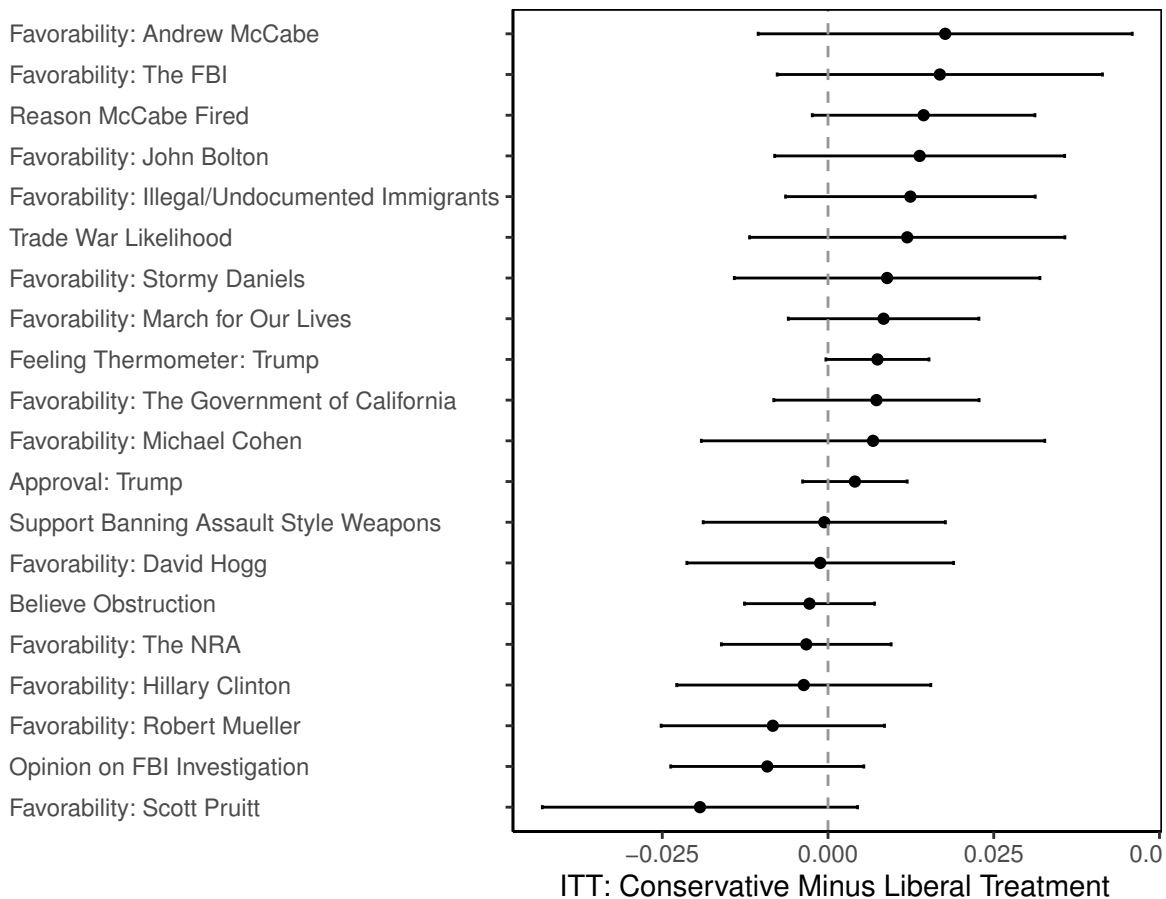
This figure shows the effect of the liberal and conservative treatments on the mean slant, in standard deviations, of all news participants engaged with, excluding the four potential liberal outlets and the four potential conservative outlets defined for each participant. The regressions control for the outcome in baseline if it exists. The sample includes 1,699 participants who installed the extension and provided access to their shared posts for at least two weeks following the intervention. Error bars reflect 90 percent confidence intervals.

Figure A.9: Effects of the Conservative Treatment on Mean Slant by Month, Compared to Liberal Treatment



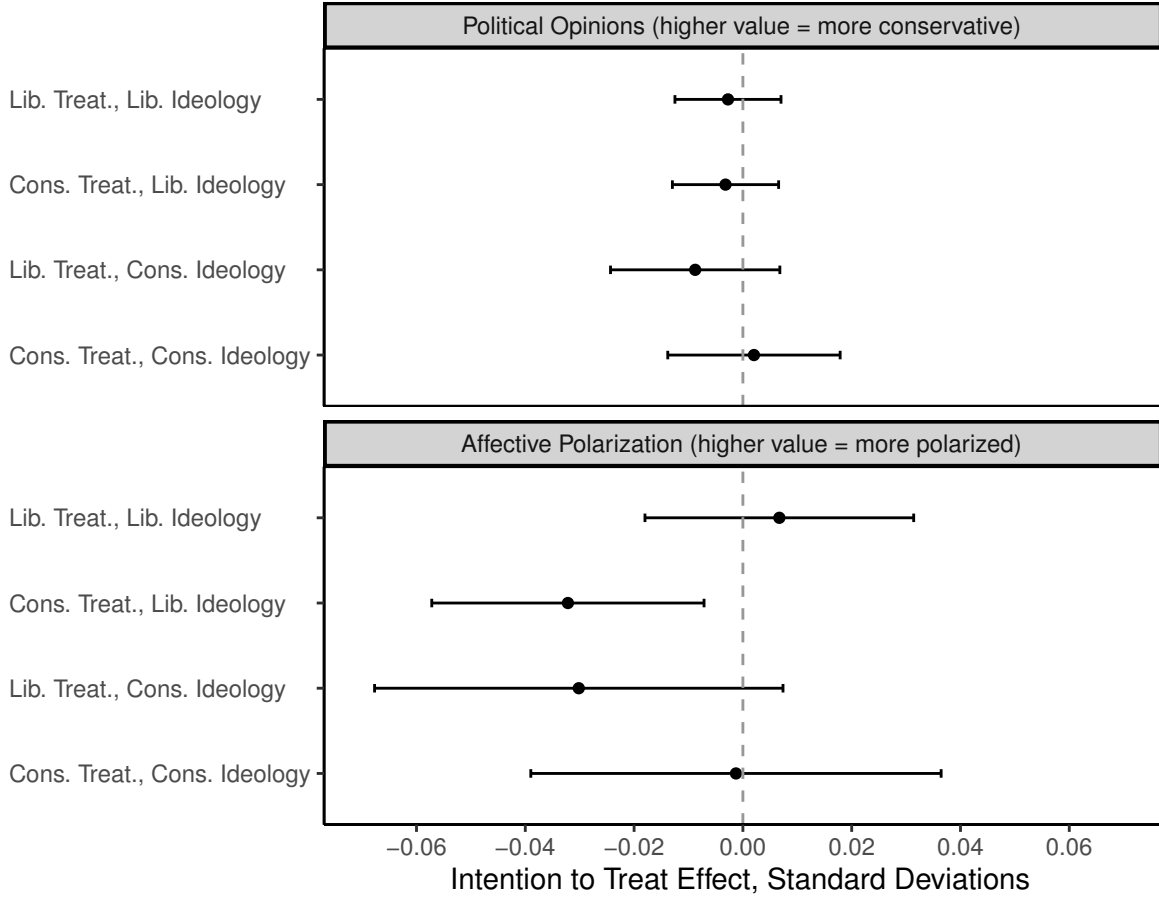
These figures show the difference between the effect of the liberal and conservative treatments on the mean slant over time. Each panel presents a series of regressions, where the dependent variable is the slant of outlets in a specific month. In the x-axis, relative month 1 is defined as 28 days immediately following the intervention. Sub-figure (a) is based on 1,351 participants who kept the extension installed for at least 84 days following the intervention. Sub-figure (b) is based on 9,932 participants who provided access to posts they shared for at least 84 days following the intervention. The regressions control for the outcome in baseline, if it exists. Error bars reflect 90 percent confidence intervals.

Figure A.10: Effects on Components of the Political Opinion Index



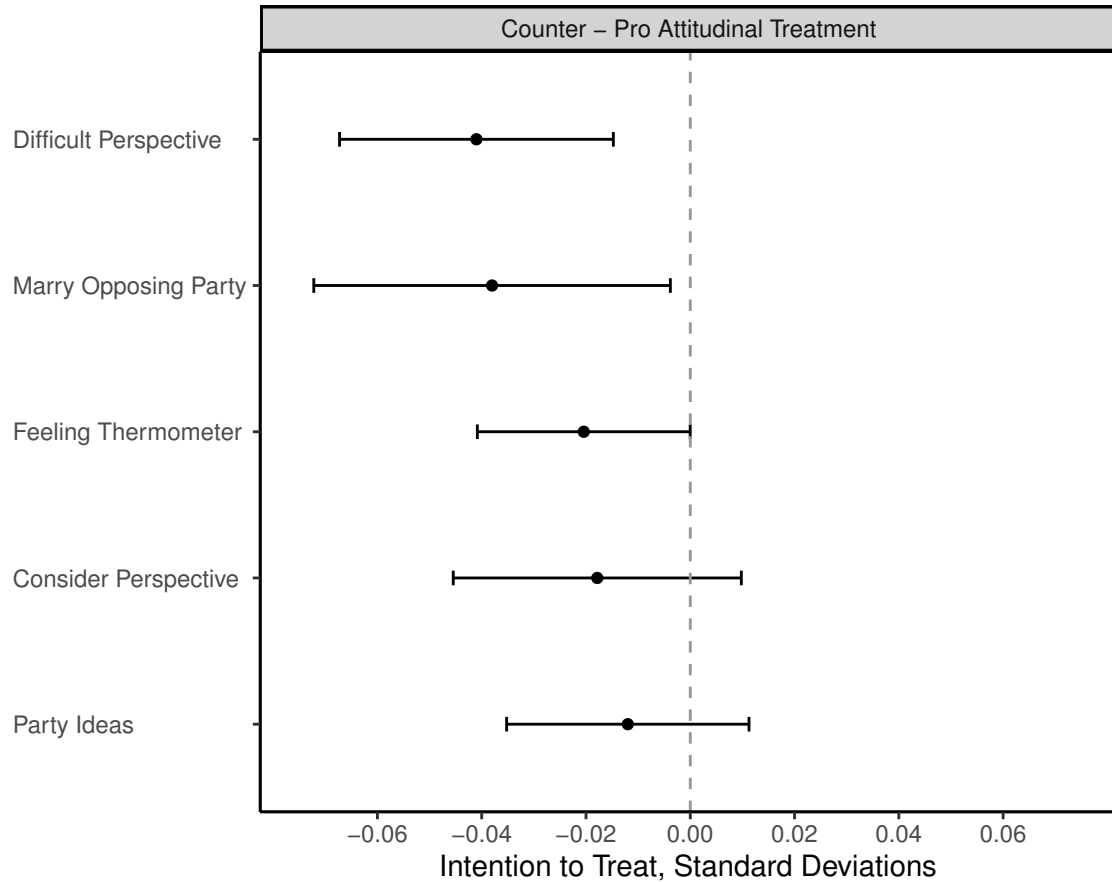
This figure shows the effect of the conservative treatment, compared to the liberal treatment on outcomes composing the political opinions index. Each row represents a separate regression as specified in Section II.E. Outcomes are defined such that a higher value is associated with a more conservative opinion and then standardized with respect to the control group. *Favorability* outcomes are based on questions asking participants whether they have a very favorable, favorable, unfavorable, or very unfavorable opinion on specific individuals or organizations. *Approval: Trump* is whether participants strongly approve, somewhat approve, somewhat disapprove, or strongly disapprove of the job Donald Trump is doing as President. *Feeling Thermometer: Trump* is feeling toward Trump on a 0-100 degrees scale. *Believe Obstruction* is whether participants believed that President Trump has attempted to derail or obstruct the investigation into the Russian interference in the 2016 election. *Opinion on FBI Investigation* is whether participants think the FBI investigation into Trump campaign officials' contacts with Russian government officials is a serious attempt to find out what really happened, a politically-motivated attempt to embarrass Donald Trump or equally-motivated by both of these. *Reason McCabe Fired* is whether participants believe McCabe was fired because of improper actions while serving as Deputy Director of the FBI, as a way to damage McCabe's credibility in any evidence he might give to the Robert Mueller investigation, or as an act of revenge (multiple choice question). *Trade War Likelihood* is whether participants believe it is very likely, somewhat likely, somewhat unlikely, or very unlikely that a trade war will develop between the United States and foreign countries in the next year. *Support Banning Assault Style Weapons* is whether participants strongly support, support, oppose, or strongly oppose banning assault-style weapons. Error bars reflect 90 percent confidence intervals.

Figure A.11: Effect of the Treatments on Primary Outcomes, by Ideological Leaning



This figure shows the effect of the interaction of treatment and ideological leaning on the primary outcomes: $Y_i = \beta_1 T_i^L I_i^L + \beta_2 T_i^C I_i^L + \beta_3 T_i^L I_i^C + \beta_4 T_i^C I_i^C + \alpha X_i + \varepsilon_i$ where: T_i^C, T_i^L are binary indicators for the conservative and liberal treatments and I_i^C, I_i^L are binary indicators for whether the participant's ideological leaning is conservative or liberal. The reference group is the control group. The controls and the definition of ideological leaning are specified in Section II.E. In the first panel, the x-axis is the ITT effect on the political opinions index, where a higher value is a more conservative outcome. In the second panel, the x-axis is the ITT effect on the affective polarization index, where a higher value is a more polarized outcome. Error bars reflect 90 percent confidence intervals.

Figure A.12: Effect of the Treatments on Components of the Affective Polarization Index



This figure shows the effect of the counter-attitudinal treatment on the measures composing the affective polarization index, compared to the pro-attitudinal treatment. Each row presents the result of a regression estimating the effect of the treatment on one dependent variable where a higher value is associated with a more polarized outcome. *Difficult Perspective* and *Consider Perspective* measure political empathy. The former is the difference in how difficult it is to see things from each party's point of view, and the latter is the difference in how important it is to consider the perspective of each party. *Marry Opposing Party* is how participants would feel if their son/daughter married someone from the opposing party. *Feeling Thermometer* is the difference in how warm participants feel toward each party. *Party Ideas* is the difference in how many good ideas each party is perceived to have. The outcomes are described in more detail in Section II.D.2 and the regressions are specified in Section II.E. Error bars reflect 90 percent confidence intervals.

Figure A.13: Recruitment Ads

(a) Political Ad

Yale Media Survey
Sponsored (demo) · 🌐

Participate in a short Yale University research survey and you can win an \$80 Amazon gift card



Interested in Politics?
Share your opinion!
YALESURVEY.QUALTRICS.COM [Learn More](#)

👍🤔👎 103 87 Comments 38 Shares

👍 Like 💬 Comment ➦ Share 🧑 ▼

(b) General Ad

Yale Media Survey
Sponsored (demo) · 🌐

Participate in a short Yale University research survey and you can win an \$80 Amazon gift card

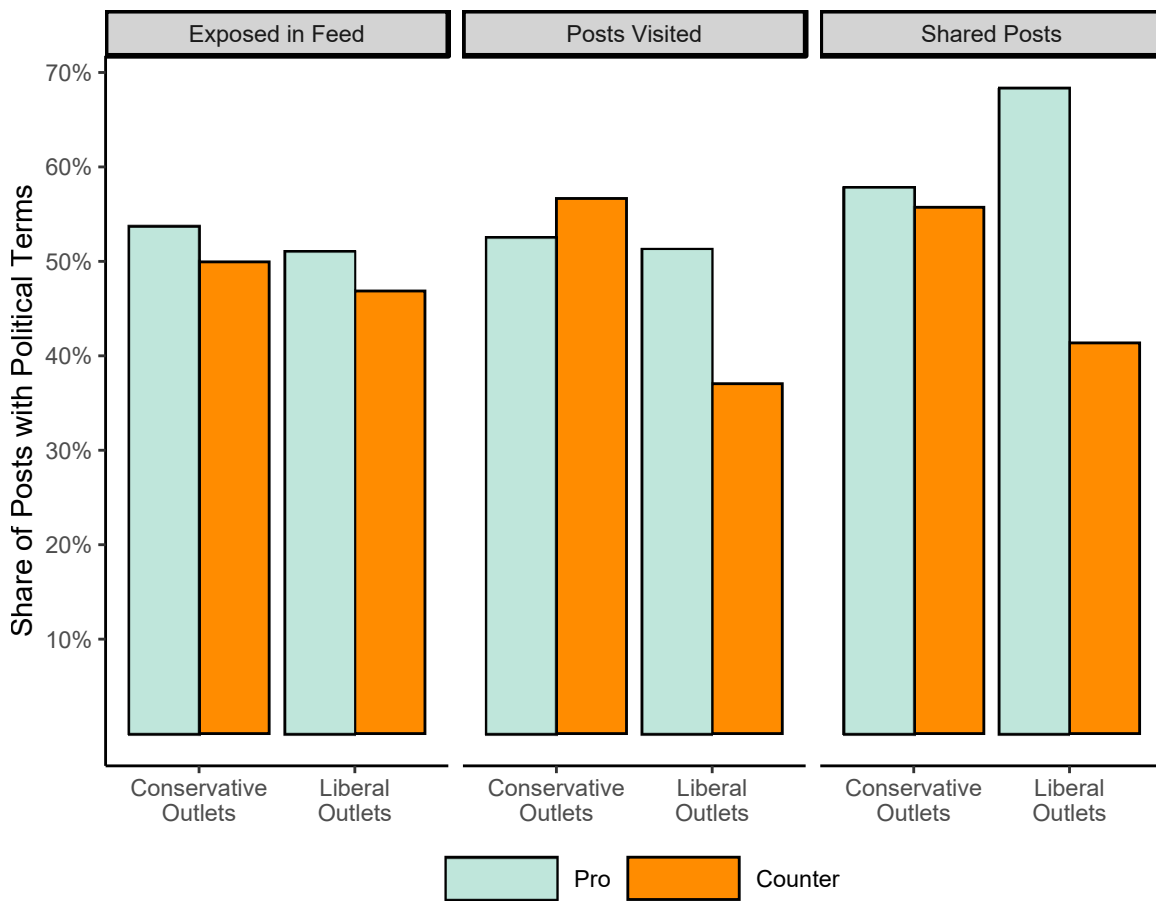


Help us understand American society better
Share your opinion and you can win an Amazon gift card!
YALESURVEY.QUALTRICS.COM [Learn More](#)

👍🤔👎 141 119 Comments 50 Shares

👍 Like 💬 Comment ➦ Share 🧑 ▼

Figure A.14: Share of Posts Mentioning Political Terms



This figure shows the share of posts mentioning political terms in posts from outlets participants subscribed to. Posts are defined as political if they contain the following terms: ar 15, biden, bolton, carson, clinton, congress, conservative, daca, democrat, devos, dnc, elect, gop, gun control, gun law, gun right, immigration, kushner, liberal, manafort, mass shooting, mccabe, mcconnell, michael cohen, nra, obama, parkland, pelosi, pence, politic, pruit, republican, sanctuary city, sanctuary state, school shooting, senate, tax cut, the left, the right, tillerson, trump, vote, walkout, white house. Posts from the pages of the eight primary outlets and first two alternative outlets (excluding suspected ads) in the first eight weeks following the intervention are included. Political terms are searched for in the post's text, URL, and any commentary included by the participants for shared posts.

Figure A.15: Links in Posts Observed in the Feed, by Outlet and Section



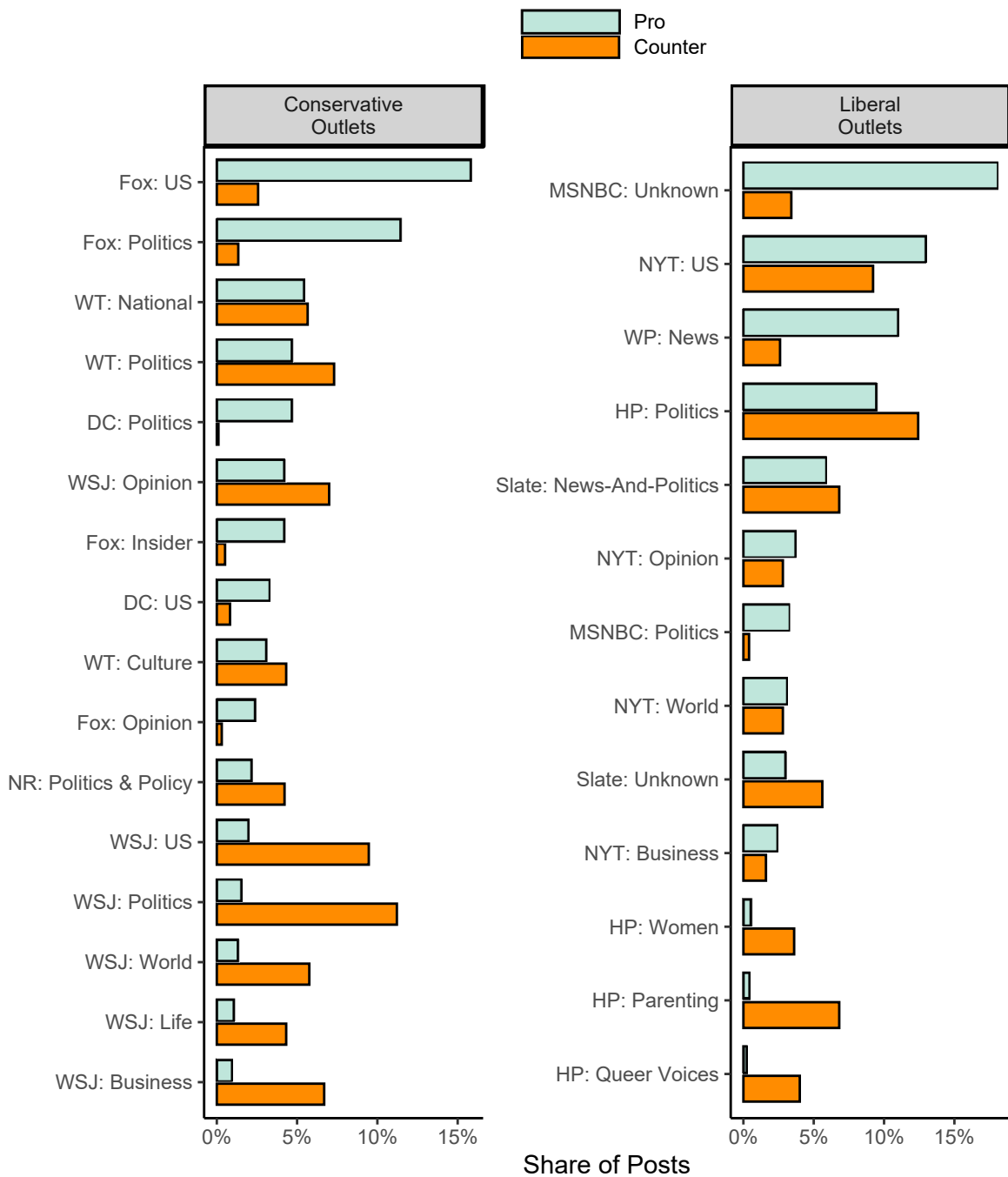
This figure shows the most common outlets and sections of links participants were exposed to in their feed. Data is from the eight weeks following the intervention. Posts from the pages of the eight primary outlets and first two alternative outlets (excluding suspected ads) are included: Daily Caller (DC), Fox News (Fox), HuffPost (HP), MSNBC, Slate, National Review (NR), New York Times (NYT), Wall Street Journal (WSJ), Washington Post (WP), and Washington Times (WT).

Figure A.16: Links Visited by Participants, by Outlet and Section



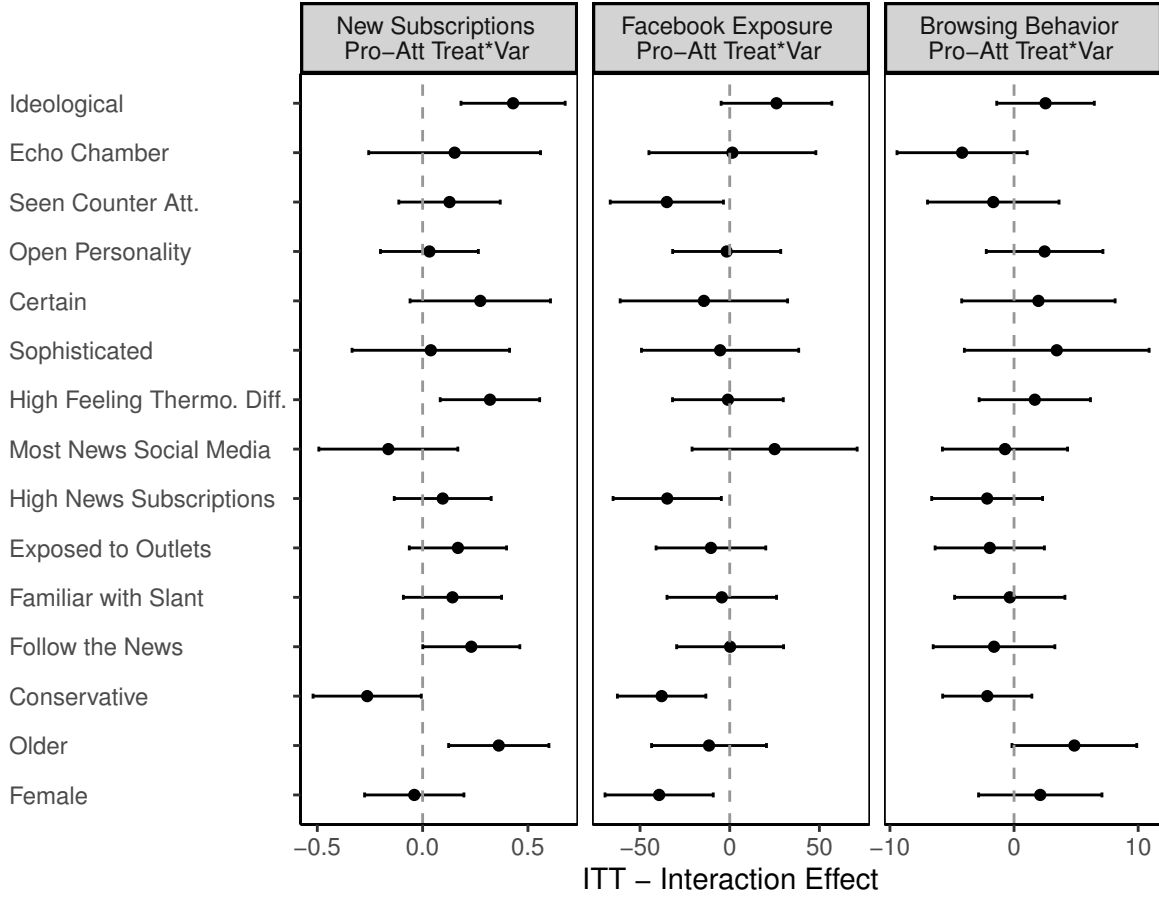
This figure shows the most common outlets and sections participants visited through links shared by the outlets they subscribed to. For more details see Figure A.15.

Figure A.17: Links in Posts Shared by Participants, by Outlet and Section



This figure shows the most common outlets and sections of the links participants shared when sharing posts from the outlets they subscribed to. For more details see Figure A.15.

Figure A.18: Heterogeneous Effects on Engagement with Pro-Attitudinal Outlets

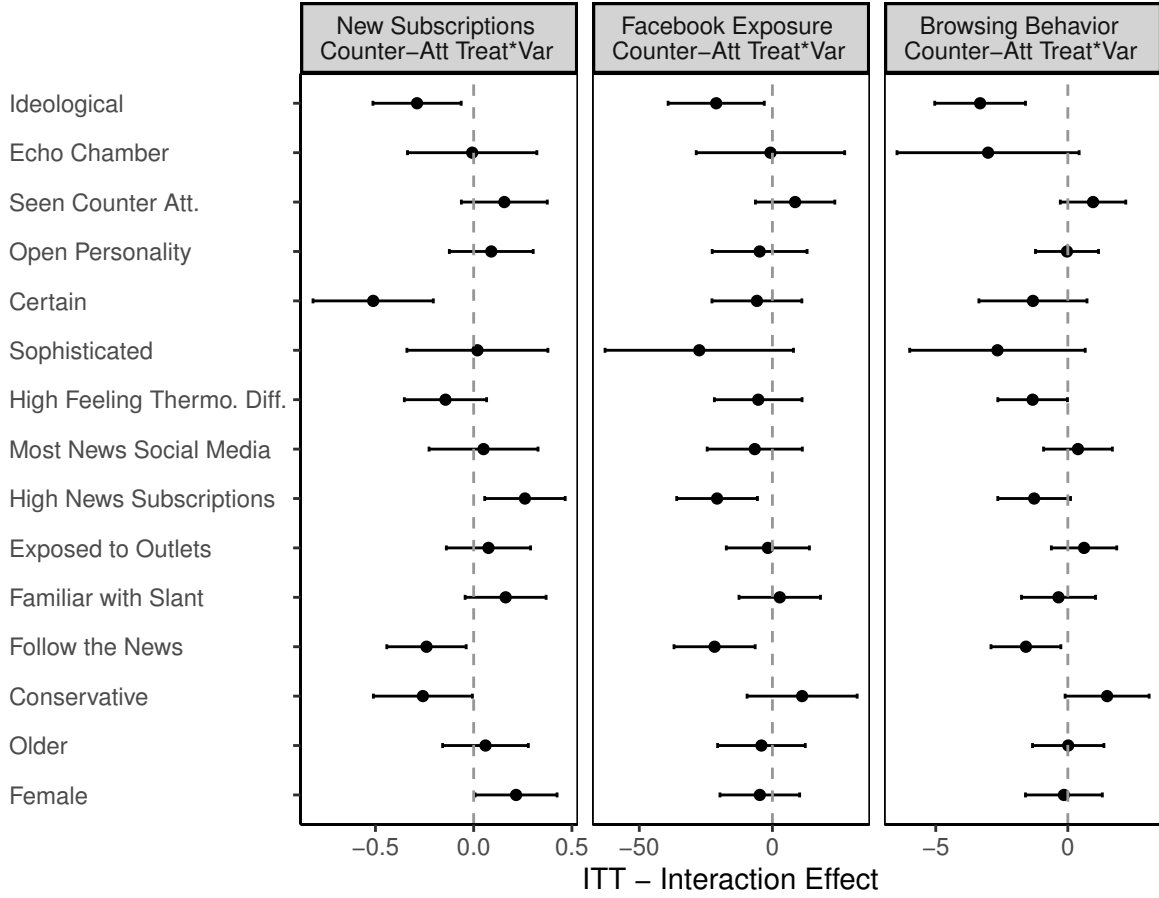


This figure shows heterogeneous effects of the pro-attitudinal treatment on engagement with the pro-attitudinal outlets. Each row presents the β coefficient in the following regression:

$$Y_i = \alpha T_i^P + \beta T_i^P \times Var + \gamma Var + \delta X_i + \varepsilon_i,$$

where the dependent variables are the number of potential pro-attitudinal outlets participants subscribed to (left panel), the number of posts from these outlets appearing in their feed (center panel), and the number of websites associated with these outlets they visited (right panel). The regressions control for the set of potential outlets defined for each participant and baseline outcomes if they exist. A higher value means individuals were more likely to engage with pro-attitudinal outlets as a result of the pro-attitudinal treatment, compared to the control group. The definitions of the variables analyzed are described in Section C.3. Error bars reflect 90 percent confidence intervals.

Figure A.19: Heterogeneous Effects on Engagement with Counter-Attitudinal Outlets

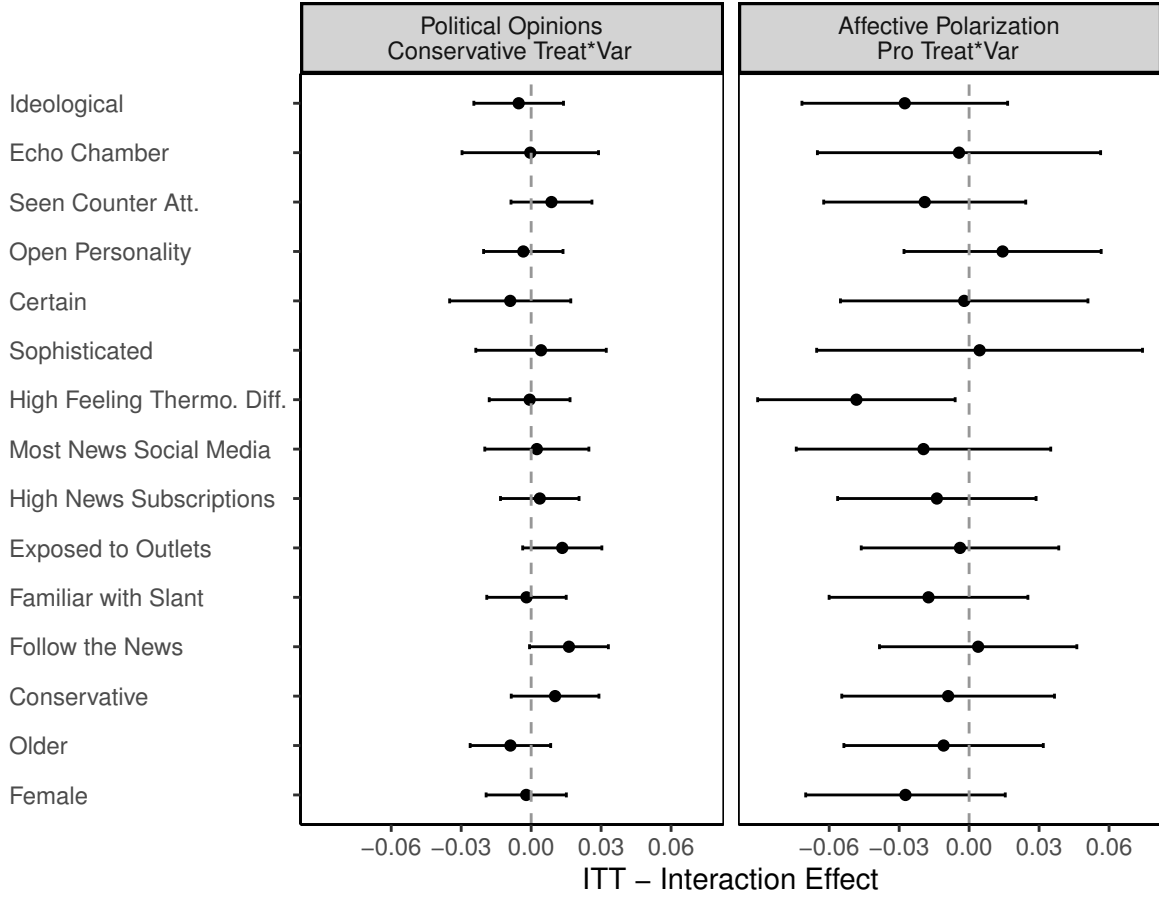


This figure shows heterogeneous effects of the counter-attitudinal treatment on engagement with the counter-attitudinal outlets. Each row presents the β coefficient in the following regression:

$$Y_i = \alpha T_i^A + \beta T_i^A \times Var + \gamma Var + \delta X_i + \varepsilon_i,$$

where the dependent variables are the number of potential counter-attitudinal outlets participants subscribed to (left panel), the number of posts from these outlets appearing in their feed (center panel), and the number of websites associated with these outlets they visited (right panel). The regressions control for the set of potential outlets defined for each participant and baseline outcomes if they exist. A higher value means individuals were more likely to engage with counter-attitudinal outlets as a result of the counter-attitudinal treatment, compared to the control group. The definitions of the variables analyzed are described in Section C.3. Error bars reflect 90 percent confidence intervals.

Figure A.20: Heterogeneous Effects on Political Opinions and Affective Polarization



This figure shows heterogeneous effects on political opinions and affective polarization. In the left panel, each row represents the β coefficient in the following regression:

$$Y_i = \alpha T_i^C + \beta T_i^C \times Var + \gamma Var + \delta X_i + \varepsilon_i,$$

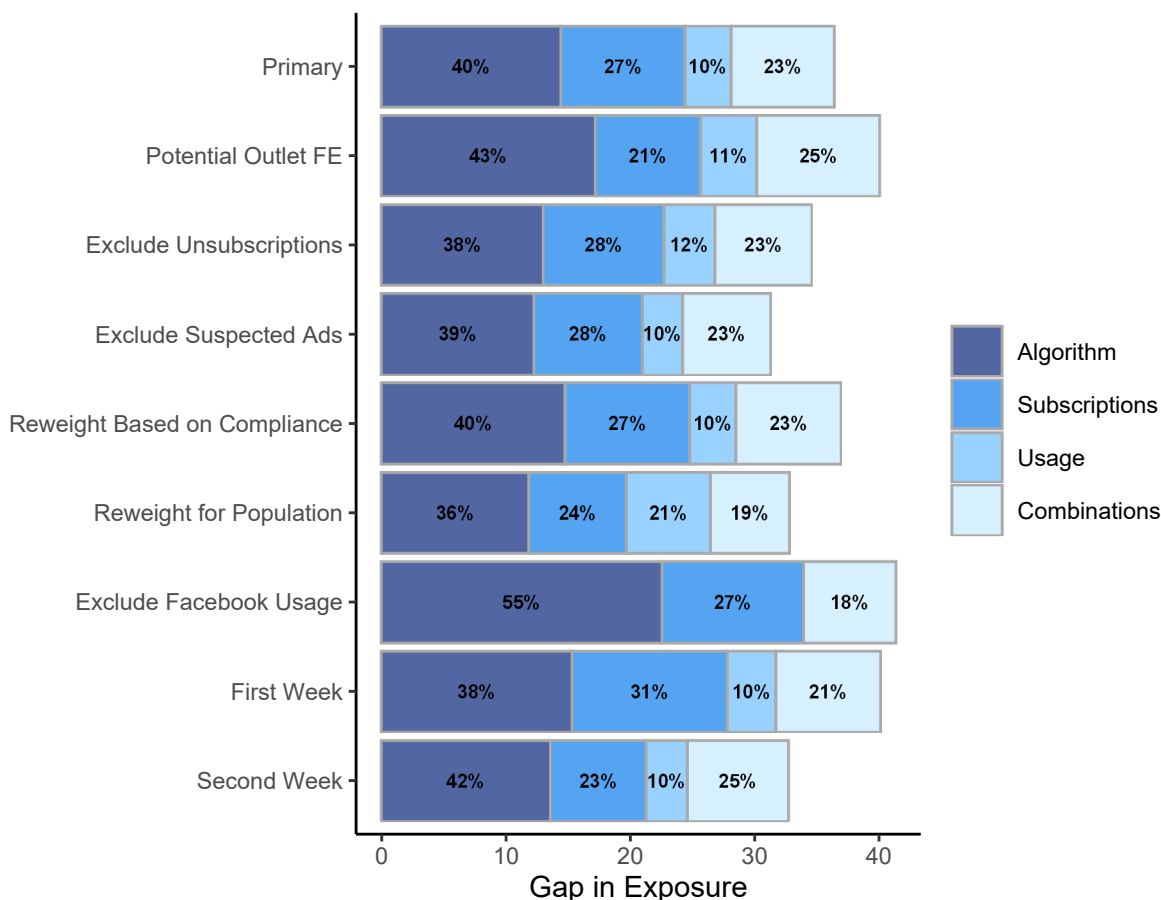
where the dependent variable is the political opinion index, and the independent variable is the full interaction of the conservative treatment and the variable analyzed in the row. A higher value means individuals were more likely to become more conservative by the conservative treatment, compared to the liberal treatment.

In the right panel, each row presents the β coefficient in the following regression:

$$Y_i = \alpha T_i^P + \beta T_i^P \times Var + \gamma Var + \delta X_i + \varepsilon_i,$$

where the dependent variable is the affective polarization index, and the independent variable is the full interaction of the pro-attitudinal treatment and the variable analyzed in the row. A higher value means individuals were more likely to become polarized as a result of pro-attitudinal treatment, compared to the counter-attitudinal treatment. The regressions control for the covariates specified in Section II.E along with the potential outlets defined for each participant. The definitions of the variables analyzed are described in Section C.3. Error bars reflect 90 percent confidence intervals.

Figure A.21: Decomposing the Gap Between Exposure to Posts from the Offered Pro-attitudinal and Counter-attitudinal Outlets, Additional Estimations



This figure decomposes the gap between the number of posts participants were exposed to from the offered pro- and counter-attitudinal outlets. The first row repeats the main specification described in Figure 10. The second row controls for the potential outlets defined for each participant. The third row defines subscriptions as subscribing to the outlet for at least two weeks. The fourth row excludes posts that are likely to be sponsored (ads). The fifth row reweights the participants in each treatment such that the compliers resemble the entire sample. The sixth row reweights the participants such that the entire sample resembles the US population. The seventh row excludes differences in usage between the groups. The final two rows decompose the results separately for the first and second week after the intervention. Each row is described in more detail in Section C.7.2.

Table A.1: Outlets Offered

Outlet	Treatment	Slant	Potential	Offered	Sub.
The Washington Times	Conservative	0.70	37,120	12,366	3,278
The National Review	Conservative	0.90	36,168	12,057	2,953
The Wall Street Journal	Conservative	0.28	35,406	11,805	4,059
Fox News	Conservative	0.78	32,566	10,842	1,425
The Daily Caller	Conservative	0.87	4,522	1,471	323
Washington Examiner	Conservative	0.82	1,719	607	133
The Western Journal	Conservative	0.90	1,531	509	153
Townhall	Conservative	0.93	397	135	37
The Blaze	Conservative	0.89	221	80	25
The Conservative Tribune	Conservative	0.89	204	72	34
Newsmax	Conservative	0.77	114	32	14
Slate	Liberal	-0.68	35,206	11,738	3,008
MSNBC	Liberal	-0.81	35,091	11,688	2,786
HuffPost	Liberal	-0.62	31,927	10,643	2,359
The New York Times	Liberal	-0.55	30,337	10,145	3,376
Washington Post	Liberal	-0.26	8,234	2,824	1,341
Salon	Liberal	-0.88	5,119	1,668	595
Daily Kos	Liberal	-0.90	2,015	661	232
The Atlantic	Liberal	-0.54	636	203	116
Mother Jones	Liberal	-0.87	515	150	59
NPR	Liberal	-0.61	431	119	70
The New Yorker	Liberal	-0.76	317	105	65
PBS	Liberal	-0.54	134	40	23

This table shows the list of outlets included in the experiment. *Slant* is the outlet's slant, ranging from -1 to 1 (Bakshy, Messing and Adamic, 2015). *Potential* is the number of participants for whom the outlet was defined as a potential outlet. *Offered* is the number of participants who were offered to subscribe to the outlet, based on their treatment assignment. *Sub.* is the number of participants who subscribed to each outlet in the intervention. The first four liberal outlets and the first four conservative outlets are the primary outlets offered in the experiment and the rest of the outlets are the alternative outlets offered if a participant already subscribed to a primary outlet. Data is from the baseline sample.

Table A.2: Descriptive Statistics by Sample

		Baseline Sample	Access Posts Subsample	Endline Survey Subsample	Extension Subsample
1)	Ideology (-3, 3)	-0.61	-0.61	-0.71	-0.95
2)	Ideology, Abs. Value (0, 3)	1.75	1.75	1.80	1.81
3)	Democrat	0.38	0.38	0.40	0.44
4)	Republican	0.17	0.17	0.16	0.14
5)	Independent	0.37	0.36	0.36	0.36
6)	Feeling Therm., Difference	50.22	50.27	50.32	51.08
7)	Difficult Pers., Difference	1.92	1.92	1.96	1.92
8)	Most News Social Media	0.18	0.18	0.17	0.16
9)	Took Survey Mobile	0.67	0.67	0.63	0.00
10)	Female	0.52	0.52	0.52	0.49
11)	Age	47.69	47.65	48.78	52.47
12)	Total Subscriptions	474	474	472	481
13)	News Outlets Subscriptions	8.11	8.11	8.28	8.61
14)	Compliance	0.53	0.53	0.58	0.76
15)	N	37,494	34,592	17,635	1,835

This table presents descriptive statistics by subsample. *Baseline Sample* includes all participants. *Access-Posts Subsample* includes participants who provided access to posts they shared for at least two weeks. *Endline Survey Subsample* includes participants who completed the endline survey. *Extension Subsample* includes participants who installed the browser extension for at least two weeks. *Ideology, Abs. Value* is the absolute value of self-reported ideology. *Feeling Therm., Difference* is the difference between feelings toward the participants' party and the opposing party according to the feeling thermometer questions. *Difficult Pers., Difference* is the difference in whether participants find it difficult to see things from the point of view of the opposing party and their own party. For all other variables, see Table 2.

Table A.3: Balance Table, Pro- and Counter-Attitudinal Treatments

Variable	Mean		Difference		
	Sample N=36,330	US	Control - Pro.	Control - Counter.	Pro. - Counter.
Baseline Survey					
Ideology, Abs. Value (0, 3)	1.80	1.31	0.00	-0.00	-0.00
Democrat	0.39	0.37	0.01	0.00	-0.01
Republican	0.17	0.30	0.00	-0.01	-0.01
Independent	0.36	0.29	-0.01*	0.00	0.01**
Vote Support Clinton	0.54		-0.00	-0.00	0.00
Vote Support Trump	0.27		0.00	0.00	0.00
Feeling Therm., Difference	50.22	38.44	0.36	0.41	0.05
Difficult Pers., Difference	1.92		0.03	0.02	-0.02
Facebook Echo Chamber	1.20		0.00	-0.01	-0.01
Follows News	3.36	2.48	0.01	0.01	0.01
Most News Social Media	0.17	0.12	0.00	-0.00	-0.01
Device					
Took Survey Mobile	0.67		-0.01*	-0.00	0.01*
Facebook					
Female	0.52	0.52	-0.01	-0.00	0.00
Age	47.91	47.70	0.02	0.08	0.06
Total Subscriptions	473		6.91	3.16	-3.75
News Outlets Slant, Abs. Value	0.54		-0.00	-0.00	0.00
Access Posts, Pre-Treat.	0.98		0.00	0.00	-0.00
Attrition					
Took Followup Survey	0.47		0.03***	0.03***	0.00
Access Posts, 2 Weeks	0.92		0.01	0.00	-0.00
Extension Install, 2 Weeks	0.05		0.00	-0.00	-0.00
F-Test			1.23	0.80	0.99
P-value			[0.20]	[0.75]	[0.48]

This table presents descriptive statistics by whether participants were assigned to the pro-attitudinal treatment, counter-attitudinal treatment, or control group. The second column shows summary statistics for American adults for whom an ideological leaning can be defined. *Ideology, Abs. Value* is the absolute value of self-reported ideology. *Feeling Therm., Difference* is the difference between the feeling toward the participants' party and the opposing party. *Difficult Pers., Difference* is the difference in whether participants find it difficult to see things from the point of view of the opposing party and their own party. *News Outlets Slant, Abs. Value* is the absolute value of the mean slant of all outlets participants subscribed to on Facebook in baseline, where slant ranges from -1 to 1. For all other variables see Table 2. Data sources for the US are specified in Appendix C.4.1. *p<0.1 **p<0.05 ***p<0.01

Table A.4: Balance Table, Liberal and Conservative Treatments, Among Participants Who Completed the Follow-up Survey

Variable	Mean			Difference		
	Sample N=17,635	US	FB Users	Control - Lib.	Control - Cons.	Cons. - Lib.
Baseline Survey						
Ideology (-3, 3)	-0.71	0.17		-0.01	-0.02	0.01
Democrat	0.40	0.35	0.30	0.01	0.01	0.01
Republican	0.16	0.28	0.21	0.00	0.00	0.00
Independent	0.36	0.32	0.35	-0.02*	-0.01	-0.01
Vote Support Clinton	0.55			-0.00	-0.00	-0.00
Vote Support Trump	0.25			0.01	-0.00	0.01
Feeling Therm., Rep.	27.54	43.06		0.20	-0.04	0.24
Feeling Therm., Dem.	47.79	48.70		0.43	0.68	-0.25
Difficult Pers., Rep. (1, 5)	3.18			0.04	0.01	0.04
Difficult Pers., Dem. (1, 5)	2.35			-0.01	-0.03	0.03
Facebook Echo Chamber	1.20		1.12	0.01	-0.01	0.01
Follows News	3.38	2.42		0.02	0.02	-0.00
Most News Social Media	0.17	0.13		-0.01**	-0.00	-0.01*
Device						
Took Survey Mobile	0.63			-0.01	0.01	-0.01
Facebook						
Female	0.52	0.52	0.55	-0.01	-0.00	-0.00
Age	48.78	47.30	42.86	0.55*	-0.31	0.86**
Total Subscriptions	472			2.37	15.27	-12.90
News Outlets Slant (-1, 1)	-0.20			0.00	-0.01	0.01
Access Posts, Pre-Treat.	0.98			0.00	0.00*	-0.00
F-Test				1.15	0.97	1.32
P-Value				[0.29]	[0.49]	[0.16]

This table presents descriptive statistics by whether participants were assigned to the liberal treatment, conservative treatment, or control group among participants who completed the endline survey. The variables are explained in the notes for Table 2. *p<0.1 **p<0.05 ***p<0.01

Table A.5: Balance Table, Pro- and Counter-Attitudinal Treatment, Among Participants Who Completed the Follow-up Survey

Variable	Mean		Difference		
	Sample N=17,130	US	Control - Pro.	Control - Counter.	Pro. - Counter.
Baseline Survey					
Ideology, Abs. Value (0, 3)	1.84	1.31	-0.00	0.00	0.00
Democrat	0.41	0.37	0.02*	0.01	-0.01
Republican	0.16	0.30	0.00	0.00	-0.00
Independent	0.35	0.29	-0.02**	-0.00	0.01
Vote Support Clinton	0.57		-0.00	0.00	0.00
Vote Support Trump	0.25		0.00	0.01	0.01
Feeling Therm., Difference	50.32	38.44	0.96*	1.10**	0.14
Difficult Pers., Difference	1.96		0.05*	0.04	-0.01
Facebook Echo Chamber	1.22		0.00	0.00	-0.00
Follows News	3.39	2.48	0.02	0.03*	0.00
Most News Social Media	0.17	0.12	-0.00	-0.01	-0.00
Device					
Took Survey Mobile	0.63		-0.01	0.01	0.01
Facebook					
Female	0.52	0.52	-0.01	-0.01	0.00
Age	48.96	47.70	0.12	0.20	0.08
Total Subscriptions	471		4.99	3.30	-1.69
News Outlets Slant, Abs. Value	0.55		-0.00	0.00	0.00
Access Posts, Pre-Treat.	0.98		-0.00	0.00	0.00
F-Test			0.63	0.75	0.57
P-value			[0.89]	[0.78]	[0.94]

This table presents descriptive statistics by whether participants were assigned to the pro-attitudinal treatment, counter-attitudinal treatment, or control group among participants who completed the endline survey. The variables are explained in the notes for Tables 2 and A.3. *p<0.1

p<0.05 *p<0.01

Table A.6: Descriptive Statistics by Compliance

	Control	All		Pro-Att.		Counter-Att.		Liberal		Conservative		
		Comply:		Comply:		Comply:		Comply:		Comply:		
		Yes	No	Yes	No	Yes	No	Yes	No	Yes	No	
1)	Ideology (-3, 3)	-0.62	-0.92	-0.27	-0.86	-0.31	-1.05	-0.25	-1.13	-0.04	-0.71	-0.51
2)	Ideology, Abs. Value (0, 3)	1.80	1.77	1.73	1.83	1.75	1.78	1.82	1.78	1.72	1.75	1.75
3)	Democrat	0.40	0.43	0.32	0.44	0.32	0.46	0.34	0.47	0.27	0.40	0.37
4)	Republican	0.17	0.13	0.21	0.15	0.21	0.12	0.23	0.11	0.25	0.16	0.18
5)	Independent	0.35	0.36	0.37	0.35	0.38	0.36	0.35	0.35	0.38	0.37	0.36
6)	Vote Support Clinton	0.54	0.60	0.44	0.60	0.46	0.64	0.46	0.65	0.39	0.55	0.50
7)	Vote Support Trump	0.27	0.20	0.34	0.23	0.34	0.17	0.36	0.15	0.38	0.25	0.29
8)	Feeling Therm., Difference	50.47	50.24	49.92	51.23	48.52	49.03	51.02	50.70	49.33	49.79	50.51
9)	Difficult Pers., Difference	1.93	1.93	1.88	1.97	1.81	1.89	1.95	1.94	1.89	1.92	1.88
10)	Facebook Echo Chamber	1.20	1.21	1.15	1.23	1.14	1.22	1.19	1.23	1.13	1.19	1.17
11)	Most News Social Media	0.17	0.18	0.17	0.17	0.17	0.19	0.17	0.18	0.17	0.17	0.17
12)	Took Survey Mobile	0.67	0.67	0.67	0.67	0.68	0.68	0.66	0.69	0.67	0.66	0.67
13)	Female	0.52	0.57	0.46	0.56	0.47	0.60	0.45	0.59	0.45	0.56	0.47
14)	Age	47.94	48.32	46.95	49.03	46.32	47.86	47.86	48.18	46.74	48.46	47.16
15)	Total Subscriptions	476	509	430	496	431	521	429	515	428	504	431
16)	News Outlets Subscriptions	8.16	8.77	7.41	8.87	7.26	8.79	7.73	8.78	7.40	8.75	7.42
17)	Certain (0, 4)	3.16	3.12	3.18	3.14	3.17	3.11	3.20	3.11	3.17	3.13	3.19
18)	Open Personality (1, 7)	5.62	5.70	5.54	5.67	5.55	5.72	5.52	5.71	5.53	5.68	5.55
19)	Seen Counter-Att. Share	0.42	0.42	0.41	0.41	0.42	0.43	0.40	0.41	0.41	0.43	0.41
20)	N	12,104	13,258	11,734	7,115	4,985	5,791	6,335	6,604	5,893	6,654	5,841

This table presents descriptive statistics on compliance by treatment arm for the entire baseline sample. *Certain* is whether participants are extremely certain (4), very certain (3), somewhat certain (2), slightly certain (1), or not at all certain (0) of their political opinions. *Open Personality* is agreement with “I see myself as open to new experiences, complex” and the reverse values of “I see myself as conventional, uncreative.” *Seen Counter-Att. Share* is the share of potential counter-attitudinal outlets the participants reported seeing in their feed among all potential outlets. The rest of the variables are explained in Table 2 and Appendix Table A.3.

Table A.7: Segregation Measures

(a) Comscore Data

Category	Share	Seg.	Slant, Abs.
1) All Browsing		0.190	0.264
2) Direct	49.9%	0.213	0.263
3) Social	5.1%	0.280	0.358
4) Search	37.3%	0.176	0.286
5) Other	7.6%	0.216	0.300
6) FB	4.2%	0.287	0.354
7) Not FB	95.8%	0.188	0.263

(b) Extension Data

Category	Share	Seg.	Slant, Abs.	Isol.	Cong.	Share Counter
1) Subscribed		0.361	0.554	0.513	0.519	0.118
2) FB Feed		0.211	0.373	0.219	0.320	0.196
3) Friends	48.2%	0.162	0.318	0.153	0.257	0.230
4) Pages	40.7%	0.283	0.449	0.366	0.398	0.153
5) Ads	11.2%	0.255	0.400	0.270	0.320	0.192
6) Browsing		0.197	0.329	0.165	0.260	0.218
7) Not FB	85.4%	0.197	0.324	0.143	0.250	0.222
8) FB	14.6%	0.222	0.361	0.252	0.308	0.203
9) Friends	53.3%	0.203	0.331	0.176	0.265	0.219
10) Pages	36.7%	0.297	0.439	0.429	0.395	0.154
11) Ads	10.0%	0.229	0.379	0.196	0.310	0.171
12) Shared		0.255	0.414	0.307	0.363	0.181

These tables display segregation measures for online and social media news engagement. Sub-table (a) is based on 2017-2018 Comscore data and sub-table (b) is based on data from control group participants in the extension subsample from the first eight weeks after the extension was installed. The segregation measures are defined in Appendix B.1. For more details on how Facebook data was processed and suspected ads were identified see Appendix A.3.

Table A.8: Additional Segregation Measures

(a) Segregation Measures Among Comscore Users Visiting News Sites Through Facebook and Through Other Sources

Category	Share	Seg.	Slant, Abs.
1) All Browsing		0.194	0.244
2) Direct	45.3%	0.217	0.252
3) Social	27.6%	0.260	0.321
4) Search	21.7%	0.147	0.252
5) Other	5.4%	0.224	0.290
6) FB	26.3%	0.264	0.325
7) Not FB	73.7%	0.186	0.236

(b) Segregation Measures Over Time, Comscore Data

Category	Share	Seg.	Slant, Abs.
1) All: 2007-2008		0.174	0.256
2) All: 2017-2018		0.190	0.264

These tables display additional measures of segregation. Sub-table (a) includes only individuals in the Comscore panel who visited multiple news sites through Facebook and through other sources. Sub-table (b) includes the 2007-2008 and 2017-2018 Comscore panels. The segregation measures are defined in Appendix B.1.

Table A.9: Segregation Measures, Visit-Level

(a) Comscore						
Category	Share	Seg.	Slant, Abs.			
1) All Browsing		0.348	0.412			
2) Direct	65.5%	0.359	0.424			
3) Social	7.3%	0.412	0.500			
4) Search	20.0%	0.264	0.352			
5) Other	7.3%	0.318	0.380			
6) FB	6.0%	0.422	0.513			
7) Not FB	94.0%	0.342	0.406			

(b) Extension Data						
Category	Share	Seg.	Slant, Abs.	Isol.	Cong.	Share Counter
1) Subscribed		0.454	0.624	0.573	0.520	0.104
2) FB Feed		0.315	0.476	0.284	0.387	0.124
3) Friends	35.8%	0.290	0.434	0.197	0.325	0.154
4) Pages	55.8%	0.331	0.504	0.458	0.428	0.107
5) Ads	8.4%	0.303	0.474	0.305	0.380	0.113
6) Browsing		0.300	0.430	0.216	0.321	0.153
7) Not FB	90.3%	0.297	0.424	0.191	0.312	0.157
8) FB	9.7%	0.323	0.485	0.373	0.405	0.113
9) Friends	43.1%	0.288	0.436	0.222	0.331	0.145
10) Pages	50.2%	0.359	0.536	0.571	0.478	0.086
11) Ads	6.6%	0.233	0.410	0.168	0.332	0.120
12) Shared		0.318	0.457	0.414	0.368	0.158

These tables display segregation measures based on visit-level data instead of aggregating data first at the user-level. In these tables users who visit more websites implicitly receive more weight. Sub-table (a) is based on 2017-2018 Comscore data and sub-table (b) is based on data from control group participants in the extension subsample from the first eight weeks after the extension was installed. The segregation measures are defined in Section III.

Table A.10: Effects of the Treatments on News Exposure, News Sites Visited and Sharing Behavior, Two Weeks Following the Intervention, Poisson Regression

	Pro-Att. Outlets Facebook Exposure (1)	Pro-Att. Outlets Browsing Behavior (2)	Pro-Att. Outlets Sharing Behavior (3)	Counter- Att. Outlets Facebook Exposure (4)	Counter- Att. Outlets Browsing Behavior (5)	Counter- Att. Outlets Sharing Behavior (6)
Pro-Att. Treat.	1.34*** (0.13)	0.29** (0.14)	0.57*** (0.21)	0.33** (0.16)	0.19 (0.25)	0.17 (0.31)
Counter-Att. Treat.	-0.06 (0.13)	-0.03 (0.14)	0.26 (0.21)	2.49*** (0.16)	0.54*** (0.19)	1.27*** (0.31)
Pro-Att. exponentiated	3.82	1.33	1.77	1.39	1.22	1.18
Counter-Att. exponentiated	0.94	0.97	1.3	12.11	1.72	3.56
Observations	1,648	1,648	1,648	1,648	1,648	1,648

This table presents the effects of the pro- and counter-attitudinal treatments on engagement with the potential pro- and counter-attitudinal outlets in the two weeks following the intervention, estimated using Poisson regressions. The sample includes participants with a liberal or conservative ideological leaning who installed the extension and provided permission to access their posts for at least two weeks following the intervention. The regressions control for the outcome measure in baseline if it exists. Robust standard error. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.11: Effect of the Treatments on News Slant by Subsample

	News Exposure			Browsing Behavior			Shared Posts		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Liberal Treatment	-0.237*** (0.060)	-0.234*** (0.063)	-0.191*** (0.073)	-0.091** (0.037)	-0.080** (0.039)	-0.100** (0.046)	-0.021* (0.012)	-0.106* (0.056)	-0.045 (0.065)
Conservative Treatment	0.355*** (0.067)	0.365*** (0.070)	0.462*** (0.082)	0.102** (0.040)	0.105** (0.041)	0.107** (0.050)	0.046*** (0.013)	0.054 (0.060)	0.131* (0.073)
Cons. Treat. - Lib. Treat.	0.59*** (0.06)	0.60*** (0.07)	0.65*** (0.08)	0.19*** (0.04)	0.19*** (0.04)	0.21*** (0.05)	0.07*** (0.01)	0.16*** (0.06)	0.18** (0.07)
Ext. Subsample	X			X			X		
Posts Subsample								X	
Ext. + Posts Subsample		X			X				X
Ext. + Posts + Endline Subsample			X			X			
Observations	1,556	1,433	1,010	1,785	1,652	1,166	18,328	979	685

This table presents the effect of the treatments on the slant of outlets participants engaged with across various subsamples. The dependent variables are the mean slant in standard deviations of posts participants were exposed to in their feed (column 1-3), of news sites they visited (columns 4-6), and of posts they shared (columns 7-9). *Ext. Subsample* refers to the extension subsample, i.e., participants who installed the extension for at least two weeks. *Posts Subsample* refers to the access posts subsample, i.e., participants who provide permissions to access their posts for at least two weeks. *Ext + Posts Subsample* refers to participants in both these subsamples. *Ext + Posts + Endline Subsample* refers to participants in these samples who also completed the endline survey. The regressions control for outcome variables in baseline when they exist. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.12: Effect of the Treatments on Primary Outcomes, Controlling for Covariates

(a) Effect on Political Opinions

	(1)	(2)	(3)	(4)
Conservative Treatment	0.010 (0.018)	-0.002 (0.006)	-0.001 (0.005)	-0.001 (0.005)
Liberal Treatment	-0.006 (0.018)	-0.009 (0.006)	-0.006 (0.005)	-0.006 (0.005)
Conservative - Lib. Treatment	0.017 (0.019)	0.007 (0.006)	0.005 (0.005)	0.005 (0.005)
Common Controls		X	X	X
Baseline Political Opinions Controls			X	X
Ex. Last Control Group Responders				X
Observations	17,635	17,635	17,635	17,237

(b) Effect on Affective Polarization

	(1)	(2)	(3)	(4)
Pro-Att. Treatment	-0.022 (0.019)	-0.003 (0.015)	0.005 (0.012)	0.005 (0.012)
Counter-Att. Treatment	-0.055*** (0.019)	-0.039** (0.015)	-0.028** (0.012)	-0.028** (0.012)
Pro-Att. Lower Lee Bound	-0.132	-0.072	-0.03	-0.012
Pro-Att. Upper Lee Bound	0.086	0.076	0.065	0.018
Counter-Att. Lower Lee Bound	-0.172	-0.115	-0.064	-0.041
Counter-Att. Upper Lee Bound	0.06	0.045	0.037	-0.016
Pro-Att. - Counter-Att. Treat	0.033* (0.019)	0.035** (0.015)	0.033*** (0.012)	0.033*** (0.012)
Common Controls		X	X	X
Baseline Polarization Controls			X	X
Ex. Last Control Group Responders				X
Observations	16,896	16,896	16,896	16,514

These tables present the effects on the political opinions and affective polarization indices. Column (1) does not control for any covariates. Column (2) controls for self-reported ideology, party affiliation, 2016 candidate supported, ideological leaning, age, age squared, and gender. Column (3), my preferred specification, also controls for baseline questions similar to endline questions composing each index. Column (4) excludes control group participants recruited to the follow-up survey with the last email sent or ad published. Without these participants, attrition is similar across treatments. To calculate Lee bounds in the specifications with control variables, I first trim the excess observation and then run the regressions with the controls. The specification and controls are described in more detail in Section II.E. Robust standard errors. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.13: Effect of the Treatments on the Affective Polarization Index, Excluding Each Index Component

	(1)	(2)	(3)	(4)	(5)	(6)
Pro-Att. Treatment	0.005 (0.012)	0.001 (0.013)	0.008 (0.013)	0.005 (0.012)	0.002 (0.013)	0.010 (0.012)
Counter-Att. Treatment	-0.028** (0.012)	-0.033** (0.013)	-0.018 (0.013)	-0.029** (0.012)	-0.035*** (0.013)	-0.020* (0.012)
Pro - Counter	0.033*** (0.012)	0.034** (0.014)	0.025** (0.013)	0.034*** (0.012)	0.038*** (0.013)	0.030** (0.012)
Excluded Measure		Feeling Thermometer	Difficult Perspective	Consider Perspective	Party Ideas	Marry Opposing Party
Observations	16,896	16,896	16,896	16,896	16,895	16,896

This table presents the effect of the treatments on the affective polarization index. Column (1) is the primary specification. In columns (2)-(6), the index is created with four of the five affective polarization index components. The specification and controls are described in more detail in Section II.E. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.14: Effect of the Treatments on Primary Outcomes, According to Outlets Offered

(a) Effect on Political Opinions, According to Outlets Offered

	(1)	(2)	(3)
Liberal Treatment	−0.006 (0.005)	−0.007 (0.005)	−0.010 (0.007)
Conservative Treatment	−0.001 (0.005)	−0.002 (0.005)	−0.007 (0.007)
Cons. Treat - Lib. Treat	0.005 (0.005)	0.005 (0.005)	0.003 (0.007)
Standard Controls	X	X	X
Potential Outlets FE		X	
Include Only Primary Outlet			X
Observations	17,635	17,635	9,630

(b) Effect on Affective Polarization, According to Outlets Offered

	(1)	(2)	(3)
Pro-Att. Treatment	0.005 (0.012)	0.004 (0.013)	−0.001 (0.016)
Counter-Att. Treatment	−0.028** (0.012)	−0.032** (0.013)	−0.031* (0.016)
Pro-Att. Treat. - Counter-Att. Treat	0.033*** (0.012)	0.036*** (0.013)	0.029* (0.017)
Standard Controls	X	X	X
Potential Outlets FE		X	
Include Only Primary Outlet			X
Observations	16,896	16,896	9,125

These tables present the effects of the treatments on the political opinions index and the affective polarization index. Column (1) is the primary specification and includes all participants. Column (2) controls for the set of eight potential liberal and conservative outlets defined for each participant. Column (3) includes only participants who did not subscribe in baseline to any of the four primary liberal outlets or the four primary conservative outlets. Thus, in this column, all participants in the liberal treatment were offered the same four primary liberal outlets and all participants in the conservative treatment were offered the same conservative outlets. The specification and controls are described in more detail in Section II.E. Robust standard errors. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.15: Effect of the Treatments on Primary Outcomes, by Subsample

(a) Effect on Political Opinions, by Subsample

	(1)	(2)	(3)	(4)
Liberal Treatment	−0.006 (0.005)	−0.007 (0.005)	−0.011 (0.018)	−0.020 (0.019)
Conservative Treatment	−0.001 (0.005)	−0.003 (0.005)	0.002 (0.018)	−0.001 (0.018)
Conservative Treat - Lib. Treat	0.005 (0.005)	0.004 (0.005)	0.013 (0.018)	0.018 (0.018)
Controls	X	X	X	X
Sample	Endline	Endline+ Posts	Endline+ Ext	Endline+ Posts+Ext
Observations	17,635	16,339	1,286	1,196

(b) Effect on Affective Polarization, by Subsample

	(1)	(2)	(3)	(4)
Pro-Att. Treatment	0.005 (0.012)	0.008 (0.013)	0.015 (0.044)	0.027 (0.046)
Counter-Att. Treatment	−0.028** (0.012)	−0.027** (0.013)	−0.072* (0.043)	−0.056 (0.045)
Pro-Att. Treat. - Counter-Att. Treat	0.033*** (0.012)	0.035*** (0.013)	0.087** (0.043)	0.083* (0.045)
Controls	X	X	X	X
Sample	Endline	Endline+ Posts	Endline+ Ext	Endline+ Posts+Ext
Observations	16,896	15,647	1,241	1,151

These tables present the effects of the treatments on the political opinions index and the affective polarization index. Column (1) is the primary specification and includes all participants who completed the endline survey. Column (2) includes only participants who also provided permissions to access their posts for at least two weeks. Column (3) includes only participants who installed the extension for at least two weeks. Column (4) includes only participants who both provided access to their posts and installed the extension. The specifications and controls are described in more detail in Section II.E. Robust standard errors. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.16: Effect of News Exposure on Affective Polarization

(a) Causal Effect Based on Experimental Variation

	IV Affective Polarization	
	(1)	(2)
FB Counter-Att. Share, Std. Dev.	-0.130* (0.067)	
FB Congruence Scale, Std. Dev.		0.105* (0.057)
Controls	X	X
First Stage F	65.1	65.22
Observations	1,072	1,072

(b) Cross-Sectional Correlation in Control Group

	OLS Affective Polarization	
	(1)	(2)
FB Counter-Att. Share, Std. Dev.	-0.385*** (0.052)	
FB Congruence Scale, Std. Dev.		0.407*** (0.054)
Data	Control Group	Control Group
Observations	352	352

These tables measure the association between exposure to pro- and counter-attitudinal news and affective polarization. *FB Counter-Att. Share* is the share of news from counter-attitudinal outlets participants were exposed to on Facebook between the baseline and endline surveys, among all news from pro- and counter-attitudinal outlets. *FB Congruence Scale* is the mean slant of all news exposed to on Facebook, multiplied by (-1) for liberal participants. Sub-table (a) shows the results of IV regressions, where the independent variables are instrumented with the treatment. The regressions control for the covariates specified in Section II.E. Sub-table (b) presents the results of regressions run only among control group participants, where the dependent variable is the affective polarization index and the independent variables are the two summary statistics (with no controls). The regressions include all participants who are both in the endline and extension subsamples and observed at least two posts from pro- or counter-attitudinal sources. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.17: Effect of the Treatments on Attitudes Toward Each Party

	Attitude Own Party (1)	Attitude Opposing Party (2)
Pro-Att. Treatment	0.008 (0.013)	-0.003 (0.014)
Counter-Att. Treatment	0.001 (0.014)	0.031** (0.014)
Pro - Counter	0.007 (0.014)	-0.035** (0.014)
Observations	16,896	16,896

This table presents the effect of the pro and counter-attitudinal treatments on attitudes toward the party the participant is associated with and the opposing party. Participants whose ideological leaning is defined as liberal are associated with the Democratic Party and participants whose ideological leaning is defined as conservative are associated with the Republican Party. The outcome for each party is an index composed of the following four questions: the feeling thermometer, how difficult it is to see things from each party's point of view, how important it is to consider the perspective of the party, and whether the party has good ideas. The controls and the definition of ideological leaning are specified in Section II.E. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.18: Primary Outcomes Using Different Index Methods

(a) Political Opinions

	(1)	(2)	(3)	(4)	(5)
Liberal Treatment	-0.006 (0.005)	-0.008 (0.017)	0.001 (0.015)	-0.007 (0.009)	-0.006 (0.007)
Conservative Treatment	-0.001 (0.005)	0.025 (0.017)	0.011 (0.014)	0.010 (0.009)	0.005 (0.007)
Cons. - Lib. Treatment	0.005 (0.005)	0.033* (0.017)	0.010 (0.010)	0.017* (0.009)	0.011 (0.007)
Controls	X	X	X	X	X
Index Method	Standard	Inv- Cov	Inv- Cov	Inv- Cov	Inv- Cov
Include Missing Outcomes	-	No	Yes	No	Yes
Replace Negative Weights With 0	-	No	No	Yes	Yes
Observations	17,635	9,434	17,635	9,434	17,635

(b) Affective Polarization

	(1)	(2)	(3)
Pro-Att. Treatment	0.005 (0.012)	0.004 (0.017)	0.001 (0.013)
Counter-Att. Treatment	-0.028** (0.012)	-0.031* (0.017)	-0.035*** (0.013)
Pro-Att. Treat. - Counter-Att. Treatment	0.033*** (0.012)	0.035** (0.017)	0.036*** (0.013)
Controls	X	X	X
Index Method	Standard	Inv- Cov	Inv- Cov
Include Missing Outcomes	-	No	Yes
Observations	16,896	10,059	16,896

These tables estimate the effects of the treatments on the primary outcomes using different summary indices. Column (1) uses equal weights for all outcomes in the index. Column (2) uses inverse-covariate weights and excludes participants with missing values for any of the index components. In Column (3), participants with missing outcomes are included with weights renormalized to sum to one, such that an outcome measure is created for all participants who have at least one non-missing outcome. Columns (4) and (5) repeat columns (2) and (3) with non-negative weights replaced with zeros and all weights renormalized to sum to one. The specifications and controls are described in Section II.E. Robust standard errors. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.19: Effect of the Treatments on Behavioral and Attitudinal Polarization Measures

	All	Affective	Behavior
Pro-Att. Treatment	0.006 (0.014)	0.005 (0.012)	−0.001 (0.018)
Counter-Att. Treatment	−0.028** (0.014)	−0.028** (0.012)	−0.010 (0.018)
Counter-Att. Treatment - Pro-Att. Treat.	0.035** (0.014)	0.033*** (0.012)	0.009 (0.019)
Controls	X	X	X
Observations	17,159	16,896	16,637

This table estimates the effects of the treatments on polarization indices. Column (1) includes the five affective components and the three behavioral components. Column (2) is the primary outcome analyzed in the paper and includes the five affective components. Column (3) includes the three behavioral components. The specification and controls are described in Section II.E. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.20: Common Phrases Mentioned When Describing the Baseline Survey's Purpose

(a) Common Three-Word Phrases by Treatment Assignment

Rank	Control	Counter	Pro
1	social media polit (0.91%)	social media polit (1.20%)	social media polit (1.36%)
2	media influenc polit (0.75%)	media influenc polit (0.94%)	media influenc polit (0.90%)
3	peopl get news (0.70%)	effect social media (0.85%)	peopl get news (0.78%)
4	peopl polit view (0.53%)	peopl get news (0.83%)	effect social media (0.66%)
5	social media influenc (0.49%)	social media influenc (0.73%)	peopl polit view (0.61%)
6	effect social media (0.46%)	social media news (0.57%)	media polit view (0.57%)
7	influenc social media (0.46%)	peopl polit view (0.56%)	social media news (0.56%)
8	media affect polit (0.44%)	media echo chamber (0.53%)	social media influenc (0.53%)
9	current polit climat (0.40%)	media polit view (0.52%)	influenc social media (0.46%)
10	social media news (0.38%)	influenc social media (0.46%)	media echo chamber (0.46%)
11	media polit view (0.38%)	media affect polit (0.41%)	polit view media (0.41%)
12	correl polit view (0.37%)	social media affect (0.40%)	social media affect (0.41%)
13	see social media (0.34%)	social media echo (0.40%)	social media effect (0.39%)
14	polit view media (0.33%)	impact social media (0.39%)	current polit climat (0.37%)
15	affect polit view (0.32%)	influenc polit view (0.38%)	influenc polit view (0.37%)

(b) Common Two-Word Phrases by Treatment Assignment

Rank	Control	Counter	Pro
1	polit view (8.31%)	social media (9.67%)	social media (9.77%)
2	social media (7.47%)	polit view (8.41%)	polit view (8.40%)
3	polit opinion (4.20%)	polit opinion (4.13%)	polit opinion (4.13%)
4	polit lean (3.39%)	news sourc (3.92%)	news sourc (3.58%)
5	news sourc (2.63%)	polit lean (3.10%)	polit lean (3.57%)
6	media polit (2.31%)	media polit (2.43%)	media polit (2.83%)
7	polit climat (1.91%)	echo chamber (2.34%)	echo chamber (1.97%)
8	polit parti (1.90%)	media influenc (1.95%)	see peopl (1.96%)
9	get news (1.69%)	see peopl (1.80%)	media influenc (1.84%)
10	media influenc (1.67%)	get news (1.74%)	media bias (1.69%)
11	media bias (1.64%)	peopl polit (1.61%)	polit parti (1.69%)
12	see peopl (1.54%)	polit parti (1.58%)	get news (1.61%)
13	liber conserv (1.47%)	polit affili (1.54%)	polit affili (1.55%)
14	peopl polit (1.45%)	polit belief (1.54%)	polit belief (1.55%)
15	polit affili (1.43%)	media bias (1.49%)	polit climat (1.55%)

These tables show phrases participants mentioned most often when asked "If you had to guess, what would you say is the primary purpose of this study?" at the end of the baseline survey. I first process the text by removing non-ascii characters, converting all characters to lowercase, removing common stop words, and stemming words to their roots. The share of responses that include the phrase appears in parenthesis.

Table A.21: Phrases with Highest Differential Usage When Describing the Survey's Purpose

(a) Control Group and the Pro-Attitudinal Treatment

Expression	Share Among Phrases with the Same Length		
	Control	Pro	Counter
chamber	0.16%	0.41%	0.47%
divers	0.01%	0.13%	0.14%
echo	0.16%	0.42%	0.47%
echo chamber	0.20%	0.51%	0.58%
media echo	0.02%	0.12%	0.13%
media echo chamber	0.02%	0.15%	0.17%
open	0.01%	0.16%	0.21%
page	0.00%	0.14%	0.19%
social	1.68%	2.21%	2.08%
social media	1.91%	2.56%	2.40%

(b) Control Group and the Counter-Attitudinal Treatment

chamber	0.16%	0.41%	0.47%
divers	0.01%	0.13%	0.14%
echo	0.16%	0.42%	0.47%
echo chamber	0.20%	0.51%	0.58%
like	0.18%	0.31%	0.46%
open	0.01%	0.16%	0.21%
page	0.00%	0.14%	0.19%
percept	0.86%	0.61%	0.50%
promot	0.03%	0.09%	0.15%
willing	0.01%	0.05%	0.10%

(c) Pro-Attitudinal Treatment and Counter-Attitudinal Treatment

connect polit	0.04%	0.07%	0.02%
like	0.18%	0.31%	0.46%
peopl identifi	0.02%	0.04%	0.01%
percept media polit	0.03%	0.04%	0
polit	10.62%	10.41%	9.67%
push	0.03%	0.07%	0.14%
push liber	0.02%	0.03%	0.09%
rang	0.02%	0.01%	0.04%
seem like	0.01%	0	0.03%
social media bias	0.03%	0.07%	0.01%

These tables show the phrases with 1, 2, 3, or 4 words with the highest differential usage between treatment arms. Differential usage is calculated using the following formula: $\chi^2 = \frac{(f_1 f_{-2} * f_2 f_{-1})^2}{(f_1 + f_2)(f_1 + f_{-1})(f_2 + f_{-2})(f_{-1} + f_{-2})}$ where f_1, f_2 are the occurrence of the phrase in the first and second groups, and f_{-1}, f_{-2} are the occurrence of all other phrases in the first and second groups. I first process the text by removing non-ascii characters, converting all characters to lowercase, removing common stop words and stemming words to their roots.

Table A.22: Most Common Two-Words Phrases Appearing in Posts

(a) Post Participants were Exposed to in their Feed

Exposed in Feed, Conservative Outlets		Exposed in Feed, Liberal Outlets	
Pro	Counter	Pro	Counter
donald trump (10.68%)	presid trump (5.15%)	presid trump (8.33%)	presid trump (7.56%)
presid donald (8.97%)	donald trump (5.09%)	donald trump (4.07%)	donald trump (4.86%)
presid trump (3.79%)	presid donald (2.92%)	white hous (3.20%)	white hous (2.66%)
white hous (2.92%)	white hous (2.58%)	stormi daniel (1.93%)	presid donald (2.16%)
high school (2.30%)	high school (1.56%)	presid donald (1.63%)	stormi daniel (2.16%)
hillari clinton (1.56%)	trump administr (1.44%)	high school (1.14%)	michael cohen (1.23%)
gun control (1.53%)	gun control (1.19%)	special counsel (1.02%)	high school (1.20%)
school shoot (1.39%)	school shoot (1.05%)	unit state (1.01%)	unit state (0.99%)
trump administr (1.33%)	special counsel (0.91%)	michael cohen (0.98%)	special counsel (0.95%)
attorney general (1.22%)	hillari clinton (0.85%)	school shoot (0.97%)	gun violenc (0.91%)

(b) Post With Links Visited by Participants

Posts Visited, Conservative Outlets		Posts Visited, Liberal Outlets	
presid trump (5.07%)	presid trump (5.33%)	presid trump (5.19%)	donald trump (3.01%)
donald trump (4.06%)	donald trump (3.18%)	donald trump (4.35%)	presid trump (3.01%)
white hous (2.84%)	white hous (3.18%)	white hous (2.12%)	day befor (0.90%)
presid donald (2.03%)	gun control (2.05%)	high school (1.17%)	former fbi (0.90%)
high school (1.83%)	hillari clinton (1.74%)	presid donald (1.06%)	high school (0.90%)
gun control (1.62%)	second amend (1.54%)	school shoot (0.78%)	someon els (0.90%)
north korea (1.42%)	presid donald (1.33%)	special counsel (0.73%)	white hous (0.90%)
attorney general (1.22%)	robert mueller (1.23%)	unit state (0.73%)	anoth child (0.60%)
hillari clinton (1.22%)	special counsel (1.23%)	michael cohen (0.67%)	anyon els (0.60%)
justic depart (1.22%)	trump administr (1.13%)	robert mueller (0.67%)	black student (0.60%)

(c) Posts Shared by Participants

Shared Posts, Conservative Outlets		Shared Posts, Liberal Outlets	
donald trump (6.37%)	presid trump (4.43%)	presid trump (9.94%)	presid trump (3.93%)
presid donald (4.51%)	donald trump (4.33%)	donald trump (4.91%)	donald trump (3.59%)
high school (4.25%)	white hous (3.75%)	white hous (3.17%)	presid donald (2.05%)
illeg immigr (4.19%)	high school (2.31%)	presid donald (1.75%)	unit state (1.20%)
hillari clinton (3.21%)	gun control (2.02%)	trump administr (1.66%)	attorney general (1.03%)
presid trump (3.00%)	presid donald (1.92%)	school shoot (1.65%)	break presid (1.03%)
trump administr (2.38%)	trump administr (1.73%)	high school (1.58%)	cambridg analytica (1.03%)
gun control (2.23%)	special counsel (1.64%)	mass shoot (1.54%)	gun violenc (1.03%)
second amend (2.02%)	gun violenc (1.44%)	stormi daniel (1.54%)	high school (1.03%)
white hous (1.61%)	robert mueller (1.44%)	robert mueller (1.51%)	school shoot (1.03%)

These tables show the most common two-word phrases mentioned in posts from the outlets that participants subscribed to. Stop word, punctuation and additional media-related words are removed and the words are then stemmed. Posts from the pages of the eight primary outlets and first two alternative outlets (excluding suspected ads) in the first eight weeks following the intervention are included.

Table A.23: Effect of the Treatments on Media Outcomes, Reweighted to Match the US Population

	News Exposure		Browsing Behavior		Shared Posts	
	(1)	(2)	(3)	(4)	(5)	(6)
Liberal Treatment	-0.237*** (0.060)	-0.337*** (0.094)	-0.091** (0.037)	-0.059 (0.052)	-0.021* (0.012)	-0.011 (0.019)
Conservative Treatment	0.355*** (0.067)	0.419*** (0.099)	0.102** (0.040)	0.148** (0.067)	0.046*** (0.013)	0.067*** (0.019)
Cons. Treat. - Lib. Treat.	0.59*** (0.06)	0.76*** (0.09)	0.19*** (0.04)	0.21*** (0.07)	0.07*** (0.01)	0.08*** (0.02)
Reweighted		X		X		X
Observations	1,556	1,556	1,785	1,785	18,328	18,328

This table estimates the effect of the treatments on the slant of posts observed in the Facebook feed, websites visited and posts shared. Columns (1), (3), and (5) show the estimates in the extension or access posts subsamples using equal weights. These columns are the same as columns (1), (4), and (7) in Appendix Table A.11. Columns (2), (4), and (6) reweight the subsamples to match the US population based on the following covariates: self-reported ideology, the share of participants identifying as Democrats, Republicans, and Independents, the difference between the participants' feelings toward their party and the opposing party, age, and the share of females. This analysis is discussed in Appendix C.4. Robust standard errors. * $p < 0.1$ ** $p < 0.05$ *** $p < 0.01$

Table A.24: Effect of the Treatments on Primary Outcomes, Reweighted to Match the US Population

(a) Political Opinions		
	(1)	(2)
Liberal Treatment	−0.006 (0.005)	−0.005 (0.007)
Conservative Treatment	−0.001 (0.005)	−0.0003 (0.008)
Cons. Treat - Lib. Treat	0.005 (0.005)	0.005 (0.008)
Controls	X	X
Reweighted		X
Observations	17,635	17,635

(b) Affective Polarization		
	(1)	(2)
Pro-Att. Treatment	0.005 (0.012)	0.019 (0.020)
Counter-Att. Treatment	−0.028** (0.012)	−0.014 (0.022)
Pro-Att. Treat. - Counter-Att. Treat	0.033*** (0.012)	0.033 (0.020)
Controls	X	X
Reweighted		X
Observations	16,896	16,896

These tables estimate the effect of the treatments on the polarization and political opinions indices after reweighting the endline participants. Column (1) uses equal weights for all participants. Column (2) reweights the participants to match the US population means based on the following covariates: self-reported ideology, the share of participants identifying as Democrats, Republicans, and Independents, the difference between the participants' feelings toward their own party and the opposing party, age, and the share of females. This analysis is discussed in Appendix C.4. The specification and controls are described in Section II.E. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.25: Predicted Effect in Full Baseline Sample

Outcome	Treatment	(1) Main Effect Estimated	(2) Predicted Effect in Subsample	(3) Predicted Effect in Baseline Sample
News exposure, posts slant (std. dev.)	Conservative treatment, compared to liberal treatment	0.592	0.545	0.571
Browsing behavior, news sites slant (std. dev.)	Conservative treatment, compared to liberal treatment	0.193	0.204	0.218
Political opinions index	Conservative treatment, compared to liberal treatment	0.005	0.003	0.003
Affective polarization index	Pro-Attitudinal treatment, compared to counter-attitudinal treatment	0.033	0.026	0.027

This table predicts the main effects estimated in the paper for the entire baseline sample. Column (1) shows the main effect estimated in each subsample. These effects are shown in columns (1) and (4) of Appendix Table A.11 and column (3) of Appendix Table A.12. For columns (2) and (3), I first estimate heterogeneous effects in the endline survey and extension subsamples using causal forests with many survey and Facebook covariates as explained in Section C.5. Column (2) predicts the treatment effect within the subsample using out-of-bag prediction. Column (3) predicts the effect for the entire baseline sample.

Table A.26: Effect of the Treatments on Self-reported Familiarity and Accurate Political Knowledge Outcomes

	Heard Michael Cohen (1)	Heard Clark Shooting (2)	Heard Louis Farrakhan (3)	Heard Clinton Speech (4)	Correct Russian Influence (5)	Correct Wall Built (6)	Correct Trump Target (7)	Correct Tax Cut (8)
Liberal Treatment	-0.004 (0.006)	0.007 (0.007)	-0.004 (0.006)	0.008 (0.008)	0.002 (0.005)	0.016* (0.009)	-0.003 (0.009)	-0.001 (0.006)
Conservative Treatment	-0.002 (0.006)	0.002 (0.007)	-0.002 (0.006)	0.019** (0.008)	0.010* (0.005)	0.0001 (0.009)	-0.007 (0.009)	0.0004 (0.006)
Cons. Treat - Lib. Treat	0.00 (0.01)	-0.01 (0.01)	0.00 (0.01)	0.01 (0.01)	0.01 (0.01)	-0.02* (0.01)	-0.00 (0.01)	0.00 (0.01)
Controls	X	X	X	X	X	X	X	X
Expected Effect	Lib Treat	Lib Treat	Cons Treat	Cons Treat	Lib Treat	Lib Treat	Cons Treat	Cons Treat
Observations	17,635	17,431	17,635	17,464	16,167	13,872	12,141	15,655

This table estimates the effect of the treatments on eight knowledge outcomes. All the outcomes are binary. *Heard Michael Cohen* and *Heard Louis Farrakhan* are whether the participant did not mark “Never heard of” when asked for their favorability ratings of the individuals. *Heard Clark Shooting* is whether the participant heard that Stephon Clark was shot and killed by police officers in Sacramento. *Heard Clinton Speech* is whether the participant heard that Hillary Clinton suggested many white women voted for Trump since they took their voting cues from their husbands. *Correct Russian Influence* is believing that “the Russian government tried to influence the 2016 presidential election”. *Correct Wall Built* is not believing that “the US has recently started building a new border wall at the US-Mexico border.” *Correct Trump Target* is not believing that “President Trump is a criminal target of Robert Mueller’s investigation.” *Correct Tax Cut* is believing that “most people will receive an income tax cut, salary increase or bonus under the new tax reform law.” All regressions control for party affiliation, ideology, vote, age, age squared, gender, whether the participant follows the news, and whether the participant stated they know the name of their representative in congress. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.27: Effect of the Treatments on Exposure to Words in the Facebook Feed

	Michael Cohen (1)	Clark Shooting (2)	Louis Farrakhan (3)	Clinton Speech (4)
Liberal Treatment	2.558*** (0.820)	1.172*** (0.350)	0.161 (0.116)	0.041 (0.041)
Conservative Treatment	0.554 (0.531)	0.080 (0.260)	0.398*** (0.103)	0.077** (0.032)
Cons. Treat - Lib. Treat	-2.00** (0.81)	-1.09*** (0.31)	0.24* (0.13)	0.04 (0.04)
Controls	X	X	X	X
Expected Effect	Lib. Treat	Lib. Treat	Cons. Treat	Cons. Treat
Observations	1,730	1,730	1,730	1,730

This table estimates the effect of the treatments on topics appearing in participants' Facebook feeds. *Michael Cohen*, *Clark Shooting*, and *Louis Farrakhan* are the number of times the terms "Michael Cohen", "Stephon Clark", and "Louis Farrakhan" appeared, respectively. *Clinton Speech* is the number of times the word Clinton appeared along with the word vote and either the word India or the word husband. All regressions control for party affiliation, ideology, vote, age, age squared, gender, whether the participant follows the news, and whether the participant stated they know the name of their representative in congress. Data is from the extension subsample from the first eight weeks following the intervention. Robust standard errors. *p<0.1 **p<0.05 ***p<0.01

Table A.28: Estimations Decomposing the Segregation in News Exposure

	Subscriptions	FB Usage: Total Posts Observed	Platform Algorithm: Share of Posts
	OLS	OLS	IV
	(1)	(2)	(3)
Pro-Att. Treatment	0.505*** (0.086)	248.765* (150.666)	
Subscriptions			0.966*** (0.093)
Subscriptions * Pro-Att.			0.460*** (0.162)
Unit	Participant	Participant	Participant* Outlet Group
Baseline Controls		X	
Mean in Counter-Att. Treatment	1.535	2043.019	0.851
Observations	1,059	1,059	2,117

This table displays the regressions used to decompose the gap in exposure to posts from the offered pro- and counter-attitudinal outlets. In column (1), the dependent variable is the number of outlets the participant subscribed to. In column (2), the dependent variable is the total number of posts observed in the feed by the participant in the two weeks following the intervention. The regression controls for Facebook visits before the intervention. In column (3), the two groups of outlets and participants are pooled in an IV regression. Each observation is a participant and the group of pro- or counter-attitudinal outlets. The dependent variable is the share of posts (in percentage points) from the group of outlets that the participant was exposed to among all posts in the participant's Facebook feed and the independent variable is the full interaction of the number of outlets the participant subscribed to among this group and whether the outlets in the group are pro-attitudinal. Subscriptions are instrumented with whether this group of outlets was offered in the experiment. The first two columns use robust standard errors and in the third column standard errors are clustered at the participant level. The sample is composed of participants who were assigned to the pro- and counter-attitudinal treatments, for which the Facebook feed is observed in the two weeks following the intervention and where at least one post is observed. *p<0.1 **p<0.05 ***p<0.01

References

- Anderson, Michael L.** 2008. "Multiple Inference and Gender Differences in the Effects of Early Intervention: A Reevaluation of the Abecedarian, Perry Preschool, and Early Training Projects." *Journal of the American Statistical Association*, 103(484): 1481–1495.
- Angrist, Joshua D., and Ivan Fernandez-Val.** 2013. "ExtrapoLATE-ing: External Validity and Overidentification in the LATE Framework." In *Advances in Economics and Econometrics - Tenth World Congress.*, ed. Daron Acemoglu, Manuel Arellano and Eddie Dekel, 401–433.
- Aronow, Peter M., and Allison Carnegie.** 2013. "Beyond LATE: Estimation of the Average Treatment Effect with an Instrumental Variable." *Political Analysis*, 21(04): 492–506.
- Bakshy, Eytan, Solomon Messing, and Lada A Adamic.** 2015. "Exposure to Ideologically Diverse News and Opinion on Facebook." *Science*, 348(6239): 1130–1132.
- Chan, Jimmy, and Wing Suen.** 2008. "A Spatial Theory of News Consumption and Electoral Competition." *Review of Economic Studies*, 75(3): 699–728.
- DellaVigna, Stefano, and Ethan Kaplan.** 2007. "The Fox News Effect: Media Bias and Voting." *The Quarterly Journal of Economics*, 122(3): 1187–1234.
- Druckman, James N., and Matthew S. Levendusky.** 2019. "What Do We Measure When We Measure Affective Polarization?" *Public Opinion Quarterly*, 83(1): 114–122.
- Flaxman, Seth R, Goel Sharad, and Justin M Rao.** 2016. "Filter Bubbles, Echo Chambers, and Online News Consumption." *Public Opinion Quarterly*, 80: 298–320.
- Gosling, Samuel D., Peter J. Rentfrow, and William B. Swann.** 2003. "A Very Brief Measure of the Big-Five Personality Domains." *Journal of Research in Personality*, 37(6): 504–528.
- Hainmueller, Jens.** 2012. "Entropy Balancing for Causal Effects: a Multivariate Reweighting Method to Produce Balanced Samples in Observational Studies." *Political Analysis*, 20(1): 25–46.
- Heckman, James J., Sergio Urzua, and Edward J. Vytlacil.** 2006. "Understanding Instrumental Variables in Models With Essential Heterogeneity." *The Review of Economics and Statistics*, 88(3): 389–432.
- Jeffrey, Lewis, B. Keith Poole, Howard Rosenthal, Adam Boche, Aaron Rudkin, and Luke Sonnet.** 2020. "Voteview: Congressional Roll-Call Votes Database."
- Peterson, Erik, Goes Shared, and Shanto Iyengar.** 2019. "Partisan Selective Exposure in Online News Consumption: Evidence from the 2016 Presidential Campaign." *Political Science Research and Methods*, 1–17.

- Rogowski, Jon C., and Joseph L. Sutherland.** 2016. "How Ideology Fuels Affective Polarization." *Political Behavior*, 38(2): 485–508.
- Shane, Frederick.** 2005. "Cognitive Reflection and Decision Making." *The Journal of Economic Perspectives*, 19(4): 25–42.
- Stone, Daniel F.** 2020. "Just a Big Misunderstanding? Bias and Affective Polarization." *International Economic Review*, 61(1): 189–217.
- Suen, Wing.** 2004. "The Self-Perpetuation of Biased Beliefs." *The Economic Journal*, 114(495): 377–396.
- Wager, Stefan, and Susan Athey.** 2018. "Estimation and Inference of Heterogeneous Treatment Effects using Random Forests." *Journal of the American Statistical Association*, 113(523): 1228–1242.