

# Ideological Segregation among Online Collaborators: Evidence from Wikipedians

Shane Greenstein  
Yuan Gu  
Feng Zhu

Working Paper 17-028



# Ideological Segregation among Online Collaborators: Evidence from Wikipedians

Shane Greenstein  
Harvard Business School

Yuan Gu  
Harvard Business School

Feng Zhu  
Harvard Business School

**Working Paper 17-028**

Copyright © 2016, 2017 by Shane Greenstein, Yuan Gu, and Feng Zhu

Working papers are in draft form. This working paper is distributed for purposes of comment and discussion only. It may not be reproduced without permission of the copyright holder. Copies of working papers are available from the author.

# **Ideological Segregation among Online Collaborators: Evidence from Wikipedians\***

Shane Greenstein, Yuan Gu, Feng Zhu

## **Abstract**

Do online communities segregate into separate conversations about “contestable knowledge”? We analyze the contributors of biased and slanted content in Wikipedia articles about U.S. politics, and focus on two research questions: (1) Do contributors display tendencies to contribute to topics with similar or opposing bias and slant? (2) Do contributors learn from experience with extreme or neutral content, and does that experience change the slant and bias of their contributions over time? Despite heterogeneity in contributors and their contributions, we find an overall trend towards less segregated conversations. Contributors tend to edit articles with slants that are the opposite of their own views, and the slant from experienced contributors becomes less extreme over time. The experienced contributors with the most extreme biases decline the most. We also find some significant differences between Republicans and Democrats.

---

\* We thank Marco Iansiti, Gerald Kane, Karim Lakhani, Abhishek Nagaraj, Frank Nagle, Michael Norton, and seminar participants at the INFORMS Annual Meeting 2015 and the Conference on Open and User Innovation 2016. We thank Justin Ng and John Sheridan of HBS Research Computing Services for their research assistance. We also thank Alicia Shems and Kate Adams for her editorial assistance. We gratefully acknowledge financial support from the Division of Research of the Harvard Business School.

## 1. Introduction

The growth of virtual communities that blur the boundaries between reader and writer has upended our understanding of processes for generating and consuming online content. These communities generate numerous cooperative and confrontational behaviors. *Contested knowledge*—which we define loosely for now as topics involving subjective, unverifiable, or controversial information—complicates the creation and consumption of content for online communities. Online communities bring together participants from disparate traditions, with different methods of expression, cultural and historical foundations for their opinions, and, potentially, bases of facts; these diverse perspectives generate challenges for online communities (e.g., Arazy et al. 2011). While many studies have examined the processes by which communities resolve conflicts, there is a lack of quantitative research about the processes in the most challenging situations, such as with debates involving contested knowledge.

In an *unsegregated* conversation, the community engages people with diverse ideas and facilitates a conversation between participants with opposing views (Benkler 2006) until participants reach a consensus. In a *segregated* conversation, like-minded participants self-select into supplying content for others with similar views and read only the content from those with whom they already agree. This behavior polarizes information consumption and sharing (e.g., Mullainathan and Shleifer 2005, Sunstein 2001), creating segregated “small villages” (e.g., Gentzkow and Shapiro 2003, Van Alstyne and Brynjolfsson 2005). Segregated conversations draw our attention because they interfere with addressing the challenges of aggregating contributions when knowledge is contested.

Our study measures the micro-behavior that supports or undermines segregated conversations in the presence of contested knowledge, characterizing the tendency of distinct types of contributors to offer slanted contributions to content that may already contain slanted content. As with Greenstein and Zhu (2012, 2016), we adapt the method developed by Gentzkow and Shapiro (2010) for rating newspaper editorials to rate the bias and slant of Wikipedia’s *content*, i.e., its articles. In these ratings, *slant* denotes degree of opinion along a continuous yardstick, from extreme degrees of red (e.g., Republican) to extreme degrees of blue (e.g., Democrat), and all the shades of purple in between. *Bias* is the absolute value of this yardstick from its zero point, and thus denotes the strength of the opinion. In comparison to the literature one of our key novelties is

our measurement of contributors: we characterize the *slant* of Wikipedia's contributors, which we define as a contributor's average propensity to make editorial changes that move articles towards more red or blue slant.

That new measure enables us to analyze two key aspects of contributor micro-behavior, namely, (1) the slant of a target selected for a contribution, and (2) the evolution of a contributor's slant over time. Specifically, we first ask: Do contributors display tendencies to edit sites with similar or opposing biases and slants from their own? If contributors tend to edit articles that agree with their own slant, we label this event *Birds of a Feather*, or BOF. In contrast, if contributors tend to suggest contributions to content that opposes their own slant, we call such a process *Opposites Attract*, or OA. We then ask whether and how experiencing BOF and OA changes the slant and bias of a contributor over time. We ask: Do contributors learn from extreme or neutral content and does that experience change the slant and bias of their contributions? Together these two tendencies characterize the propensity to have (un)segregated conversations.

The setting for this investigation is the histories of contributions to articles about U.S. politics published in Wikipedia on January 16, 2011. Wikipedia offers a rich setting for investigating micro-behavior behind segregated conversations when knowledge is contested: All revisions are well documented, and plenty of debates, especially those about political topics, involve contested knowledge. We examine the latest version of 70,305 articles about U.S. political topics, which receive contributions from 2,891,877 unique contributors. As with prior research (e.g., Greenstein and Zhu, 2012, 2016), we characterize all articles for bias and slant along a numerical yardstick. In this study we develop a rating of the bias and slant of the *contributors* by measuring how much they add to the slant of an article on average. Then we characterize how that tendency towards bias and slant evolves over time. To our knowledge, this is the first study to analyze and measure the (mis)match between the slant of contribution and content and its evolution.

This study is also of independent interest for research on online segregated political discussions. Most reference information has moved online. Across all developed economies online sources have displaced other sources of information. Wikipedia is both a top-twenty site in almost every developed country, and, by far, the most popular and referenced online repository of comprehensive information in the developed world, with the English language version of

Wikipedia receiving over 8 billion page views a month at the time we collected the data for this study.<sup>1</sup> Its prominence makes the understanding of its production important in its own right.

Wikipedia has other advantages as a setting. Wikipedia has been operating since 2001, making it one of the oldest and longest continuously operated communities producing online content. That long life enables research into the evolution of micro-behavior over time, which is novel for studies of segregated conversations. Moreover, while the Wikipedia community espouses the ideal that it aspires to achieve a neutral point of view in its content, this is more of a belief about the process than a tested fact. Little is known about whether content arises from segregated or unsegregated communities and, relatedly, whether contributors have a tendency towards BOF or OA.

The findings are striking. We show that, in spite of considerable heterogeneity, contributors on Wikipedia display an overall tendency that points towards a less segregated conversation. The heterogeneity is complex and nuanced: Contributors with every possible bias and slant contribute to articles containing every other possible bias and slant. In spite of that variance, more contributors in Wikipedia exhibit a pattern of behavior consistent with OA than with BOF. For example, a slanted contributor is on average 8% more likely to edit an article with the opposite slant than one with the same slant. In other words, contributors with different political viewpoints tend to dialogue with each other during their editing of contestable knowledge.

The second finding points in the same direction: Contributors' slant does not persist. Contributors tend to demonstrate less, not more, bias over time. The largest declines are found among contributors who edit or add content to articles that have more biases. Editing articles reduces a contributor's slant, and editing more biased content makes contributors offer less biased contributions later. Together with the first finding, this tendency reduces segregated conversations.

These findings enhance the understanding of prior work (Greenstein and Zhu 2016), which finds that revisions in Wikipedia tends to lead to more neutrality in its content, but only very slowly. Past work could not focus on the contribution of segregated conversations, however, because it had not developed measures of the slant of contributors. In contrast, this study characterizes contributor heterogeneity as well as content creation from contributors, which permits analysis of the speed of adjustment for different types of slants and biases in content. That also enables a general characterization of how adjustment processes differ over time by type of

---

<sup>1</sup> See Wikimedia Report Card, with all data reported here: <https://reportcard.wmflabs.org/>, accessed January 2017.

contributor. For example, on average, our estimates suggest it takes extreme Republican content one year longer to reach neutrality than it does for extreme Democrat content. In the study we will trace this distinction to differences in the topics where Democrats and Republican contributors participate. Also, because the study focuses on micro-behavior of contributors, it lends itself to tests of alternative explanations, aiding inferences about the causes of segregated conversations. In summary, the study permits us to conclude that segregation declines over time *because* contributors have the tendency to both add to content with opposite points of view and moderate their own contributions over time.

### *1.1. Relationship to Prior Work*

The diffusion of the web reduced the costs of assembling the attention of many reviewers and contributors, making it feasible to arrange for a crowd to focus on the same topic. That does not imply it is feasible for every topic to garner useful attention from a large crowd, however. Topics vary in the type of contributors they attract, in the viewpoints of those contributors, and the type of contributions they make. With every topic the crowd faces numerous challenges aggregating the information from many contributors into text that others find useful, readable, and accessible. We build on considerable prior work (Greenstein and Zhu 2012, 2016) and focus on a setting where the challenges are greatest: Where the knowledge is contested.

Our study of segregated conversations builds on the work of many studies of ideological segregation on the Internet (e.g., Sunstein 2001; Carr 2008; Lawrence, Sides, and Farrell 2010; Gentzkow and Shapiro 2011). The concern with segregated conversation in prior work was motivated by many reasons. Segregation can facilitate radicalization of some individuals and groups (Purdy 2015).<sup>2</sup> The persistence of many segregated conversations also can prevent varying perspectives into a common view, and delay confrontation or a political discourse between contradictory facts and ideas. It also has been held responsible for discouraging interracial friendships, disconnecting different social segments, and stimulating social isolation. Prior work emphasizes different causes, such as the role of the social network structure of online communities (e.g., Ahn et. al. 2007), and the factors that facilitate information contribution in online

---

<sup>2</sup> See, for example, <http://www.vice.com/read/we-asked-an-expert-how-social-media-can-help-radicalize-terrorists> and <http://www.rand.org/randeurope/research/projects/internet-and-radicalisation.html>, accessed October 2016.

communities (e.g., Jeppesen and Frederiksen 2006; Chiu et al. 2006; Ma and Agarwal 2007, Xu and Zhang 2013). None of this prior research focuses on how contested knowledge shapes the formation of segregated conversation, as does our study.

One line of prior work assumes a single “right” answer exists and examines whether (and how) online crowds reach that right answer (Page 2007). Several of variants on this research presume the existence of a single “consensus forecast,” and examine whether contributors herd around the consensus or deliberately choose “extreme” positions to influence the consensus (Laster, Bennet and Geoum 1999; Zitzowitz 2001). This study’s approach differs in the characterization of behavior. Prior literature presumes an extrinsic motive for herding or departing from the consensus. Our study presumes contributors have intrinsic biases – i.e., desire to express their opinions – and that motivates their contributions. Our measurement strategy also differs, because this approach requires measuring the intrinsic leaning of a contributor.

We relate particularly to research about how herding behavior in social media shapes outcomes. Prior work examines online sites that aggregate ratings and whether individuals follow their predecessors in assigning a rating (Lee et al. 2015). Research has stressed the role of group thinking (e.g., Janis 1982), decreased communication cost (Rosenblat and Mobius 2004), emotional contagion (e.g., Barsade 2002), and, broadly, the occurrence of homophily in social networks (e.g., McPherson et al. 2001). We borrow from the approach that examines the interactions between content and contributor and modify it for Wikipedia. For example, prior research asks: Does a participant’s rating/assessment align with an aggregated report of prior ratings/assessments (e.g., Muchnick et al. 2013)? By comparison, we ask: Does a contributor add to content with a slant which matches their own, and how does that behavior change over time?

This study also adds to work that focuses on the behavior of segregated online conversations. Gentzkow and Shapiro (2011) focuses on online conversations about political content and other topics, while this work focuses on measuring and characterizing outcomes – namely, how segregated communities appear to be. Relatedly, Gentzkow and Shapiro (2010) starts from the premise that there are ideological tendencies that appear in the language of speakers, and it is this insight we borrow for our framework. In traditional media, it is found that ideological bias in news content affects political behavior (e.g., DellaVigna and Kaplan 2007; Stone 2009; Chiang and Knight 2011; Durante and Knight 2012). Prior work has also stressed partisanship persistence in online media (e.g., Larcinese et al. 2007) and identified its importance for ideological segregation



in media (e.g., Carr 2008; Lawrence et al. 2010; Gentzkow and Shapiro 2011), but not tested if participants change their behavior over time, as in our study. Most other work treats the sources of bias as isolated (e.g., Groseclose and Milyo 2005; Besley and Prat 2006; Reuter and Zitzewitz 2006; Bernhardt et al. 2008) and does not link them to contested knowledge and political discourse, which this study does.

This study is the first to examine segregated conversations in the communities that produce online reference information. Wikipedia is an important site due to its heavy use, as earlier noted. Due to the success of Wikipedia's ability to aggregate contributions into a neutral point of view, our findings suggest online conversation can develop mechanisms to overcome tendencies toward segregated conversation. Our findings also suggest that some behavior supporting segregated conversation does not persist.

Our findings raise as many questions as they answer about how unsegregated conversations arise. While many participants inside Wikipedia believe its processes help its online communities meet the ideals to which the site aspires, little quantitative evidence or controlled experiments either confirms or refutes this belief. Like other online communities, Wikipedia has adopted explicit rules, norms, policies (Forte et al. 2009; Jemielniak 2014; Schroeder et al. 2012), and quality assurance procedures (Stvilia et al. 2008), which appear to shape behavior. Many online communities have adopted schemes of access privileges that formally define roles in the organization (Arazy et al. 2015; Burke et al. 2008; Collier et al. 2008; Forte et al. 2012), and so has Wikipedia. These lead to a myriad of coordination mechanisms (Kittur et al. 2007a; Kittur and Kraut 2008; Kittur et al. 2007b; Schroeder and Wagner 2012), social interactions (e.g., Halfaker et al. 2011; Forte et al. 2012), and behaviors aimed at conflict resolution (Arazy et al. 2011).

Our findings suggest Wikipedia's mechanisms are working as desired in many circumstances. Our findings leave open questions about which specific mechanisms or combination of norms are primarily responsible, and which are comparable to institutions found in other settings.

## **2. Measurement and Setting**

We begin by defining terms and offering a simple model to motivate our measurement approach to this setting. As in Gentzkow and Shapiro (2010) and Greenstein and Zhu (2012, 2016), we first define the *slant* of content. This indicates which way a particular piece of content “leans.” It takes a numerical value, bounded on the interval  $[-D, R]$ ,  $D > 0$ , and  $R > 0$ . We normalize a

neutral point of view to 0. *Bias* of content is the absolute value of slant. We define the slant and bias of a contributor in an analogous fashion.

### 2.1. Simple models of Slant in a crowd

One standard model of a crowd presumes a single objective answer, and a platform aggregates contributions from the crowd. In many models the results improve with a larger sample of contributions (Page 2007). We modify this model for a setting in which two groups of contributors aspire to improve a controversial topic and do not agree on a single objective answer.

We illustrate this point by building on one of the simplest models of crowds. In this model two groups hold opinions along a line on the interval  $[-D, R]$ , where  $2 > R/D > 1/2$ .<sup>3</sup> One set of participants hold opinions between  $[0, R]$  and the other holds opinions between  $[-D, 0]$ , and they have an irreconcilable disagreement with one another (except for a tiny set who hold a “neutral” view around zero). These opinions are built on unverifiable facts and subjective information, and views do not change when confronted with one another. Define two sets of opinions as  $O_D$  and  $O_R$ .  $O_D$  includes all potential opinions on the interval  $[-D, 0]$  and  $O_R$  is on the interval  $[0, R]$ .<sup>4</sup> The online platform aggregates contributions from a subset of contributors in either or both groups. Define the *number* of contributions from those who hold opinions within  $O_R$  and  $O_D$  as  $N_R$  and  $N_D$ , respectively, and  $N = N_R + N_D$ .

In this model each contributor has an opinion,  $o_i$ , and  $i$  indexes the sequence of contributions as 1, 2, 3, ...,  $i$ , ...,  $N$ . Define Slant,  $S_i$ , as an aggregation of the contributions of opinion. For illustrative purposes, we define the function for Slant of a topic in the simplest possible way, as the mean of all contributions to that topic. The process for determining the draw of opinions then determines Slant.

Consider a model of  $o_i$  where opinion among potential contributors follows a uniform distribution. In one standard model of crowds the contributions of opinion are iid and drawn equally from any opinion between  $-D$  and  $R$ . The randomness reflects one of the features we will see in our application, in which a substantial fraction of contributions come from individuals who make one suggestion and no more.

---

<sup>3</sup> This latter assumption is for technical purposes only. It says that the most extreme representative of one view is not substantially more biased than the most extreme of the other. This is useful for guaranteeing convergence.

<sup>4</sup> It will be convenient to include a neutral opinion in both sets, though this is not an essential feature of the model.

This model generally does not lead to a neutral outcome. The law of large numbers suggests Slant approaches  $(R - D)/2$  as  $N$  becomes large. As  $N$  becomes large, the spread around the Slant also will become tight. This outcome is neutral only in the situation when  $R = D$ . Otherwise, the slant will equal some arbitrary point in the “interior,” and eventually settle into a situation with, at most, only incremental change.<sup>5</sup> In short, there is no reason to think merely drawing opinions randomly from a crowd can lead to an aggregation of opinions that is neutral.

Following the herding literature, we next consider two simple situations in which contributions react to aggregated opinions. These illustrations modify the assumption, as in Mullainathan and Shleifer (2005), where contributors prefer to contribute to articles that are consistent with their ideological beliefs. In our setting, contributors with intrinsic ideological slants have a choice over many articles to which they can contribute. The assumption can take one of two forms in the presence of contested knowledge. In one form contributors prefer to *avoid* contributing to any article that *already* disagrees with their beliefs, so they add only to those with which they *already* agree. In another form contributors prefer to add to articles that *disagree* with their views, so their contribution changes the article, *making it closer* to their beliefs. As a simple model of each will illustrate, one of these will lead to a segregated conversation and the other will lead to an unsegregated conversation.

We use  $S_i$  to denote the slant that includes all opinions up to  $o_i$ . Define a function  $f$ , that defines the relationship between contributed opinion and the prior slant, that is,  $o_i = f(S_{i-1})$ . Consider a model of segregation. In this model, contributors prefer articles that already slant away from a neutral point of view in a direction consistent with their beliefs, and  $f$  follows a rule: *If  $S_{i-1} < 0$  then  $o_i$  is drawn i.i.d. from  $O_D$ , otherwise from  $O_R$ .*

In this simple model of segregation, the sign of the slant attracts new random contributions of the same sign. The first draw determines all subsequent contributions. If the first draw is negative, then all subsequent draws are randomly drawn negative opinions between  $-D$  and zero. The law of large numbers suggests Slant approaches  $-D/2$  as  $N_D$  becomes large.<sup>6</sup> Similarly, if the first draw is positive, then Slant approaches  $R/2$  as  $N_R$  becomes large.<sup>7</sup>

---

<sup>5</sup> The distribution of around  $S_i$  will be  $(R+D)^2/(12N_i^{1/2})$ , becoming very small as  $N$  grows large.

<sup>6</sup> The distribution around  $S_i$  will be  $D^2/(12N_{D_i}^{1/2})$ .

<sup>7</sup> The distribution around  $S_i$  will be  $R^2/(12N_{R_i}^{1/2})$ .

Next consider specification for  $f$  where contributors make alterations to articles that disagree with their views so they can contribute to altering them. In this case, contributors make alterations to articles that slant away from a neutral point of view in a direction inconsistent with their slant. This is a simple model of unsegregated conversation. Here  $f$  follows a rule: *If  $S_{i-1} > 0$  then  $o_i$  is drawn i.i.d. from  $O_D$ , otherwise from  $O_R$ .*

In the model of unsegregated conversation the sign of the slant attracts new contributions of the opposite sign. If the Slant is negative (positive), the next contribution will be positive (negative). In this case it does not matter whether the first draw is negative or positive. Contributions will move the Slant towards the center in either case. As  $N$  grows large the contribution from each contribution declines, and slant settles near zero.<sup>8</sup>

While many crowd models with contested knowledge are possible, this simple model is sufficient for illustrating several features. First, neutrality cannot emerge from a model that randomly draws from opinions. Second, the model forecasts an association between a reinforcing process and segregated conversations, i.e., contributions from those with similar slant will appear to be segregated. Third, it suggests that unsegregated conversations will display a process that does not reinforce existing slant and will draw opposite opinions. Fourth, the model suggests that segregated conversations are associated with more biased outcomes than unsegregated conversations, and the latter are associated with a comparatively moderate slant near the neutral point of opinion. Finally, the model suggests that the slant only settles down in a single place after the number of suggestions reaches a large number (albeit, it is unclear from the model precisely what “large” means in practice). These observations inform our statistical analysis below. In our application below we will discuss a specific setting in which the underlying distributions are not observable, but the sequence of contributions are, as are the resulting slants.

## *2.2 The measurement of segregated conversations*

Our measurement strategy resembles Greenstein and Zhu (2012, 2016), which builds on Gentzkow and Shapiro (2010) and adapts the strategy to Wikipedia. There is a key novelty to our

---

<sup>8</sup> If the slant is negative, then the next draw is positive. If the slant is negative again, then again the draw is positive. This continues until the slant is positive. If the slant is positive, then the next draw is negative, and so on. In this way the slant draws new opinions of the opposite sign. As  $N$  grows large, the incremental contribution cannot change the result much. At most a new opinion moves the average no more than either  $R/N$  or  $-D/N$ , which becomes small as  $N$  grows. In this way the process will approach zero.

measurement strategy: we characterize the tendencies of a contributor to a specific topic – whether a contributor tends to make edits that push the topic in a blue or red direction.<sup>9</sup> Then we analyze two endogenous choices of contributors: whether to contribute to the topic with a slant that is similar or different than their own and whether to change the slant of their contribution over time.

Some shorthand will be useful for describing empirical regularities below. Birds of a feather, or BOF, arises in two ways: When a Democratic contributor edits content with a Democratic slant, or when a Republican contributor edits content with a Republican slant. Opposites Attract, or OA, arises in two different types of situations: When a Republican contributor edits Democratic content, or when a Democratic contributor edits Republican content.

If a contributor acts in ways consistent with BOF, then additional contributions will reinforce the preexisting slant. If a majority of contributors act in accordance with BOF, then segregated conversations will arise. In contrast, if a contributor acts in ways consistent with OA, additional contributions will not reinforce the existing slant, but will reduce the bias of the content.

The discussion so far presumes a contributor retains a fixed slant over his or her lifetime of contributions. A second set of questions arise in a setting with a long history of contributions. Do contributors alter their behavior after contributing to extreme or neutral content? Does experience reduce or increase the bias of their contributions? If so, by how much? These questions have not been a focus of prior research. They arise naturally in this analysis, due to the availability of information about the long-term experience of contributors with (un)segregated conversation.

Together, the two questions can flexibly identify the micro-behavior that supports tendencies towards segregated or unsegregated conversations. In one possible extreme, contributors could display BOF and not alter the slant of their contributions over time. That would reinforce segregated conversations. If, on the one hand, contributors display OA and alter their contributions over time towards more neutrality, then conversations will tend towards a less segregated conversation. It is also possible that the two micro-behaviors could work in opposite directions, which could result in segregated or unsegregated conversations. In that sense the approach does not presume anything about the underlying micro-behavior or the outcome.

---

<sup>9</sup> Our measurement strategy uses text-based keywords to measurement slant and bias. This contrasts with citation-based measures of slants, such as Groseclose and Milyo (2005). They count the times that a media outlet cites a list of 200 think tanks in the United States and then compare this with the times that members of Congress cite the same think tanks in their speeches on the floor of the House and Senate. We cannot use this method because our analysis examines individual articles on Wikipedia and most of them do not cite these think tanks.

This approach also can potentially migrate to any setting with segregated and unsegregated communities. As we describe below, BOF and OA are identified under weak and plausible assumptions about the exogeneity of existing content’s slant/bias to a contributor and under mild assumptions about a contributor’s slant/bias following standard statistical properties.

### 2.3 *Empirical setting*

Founded in 2001, Wikipedia positions itself as “the free encyclopedia that anyone can edit”—that is, as an online encyclopedia entirely written and edited via user contributions. Topics are divided into unique pages, and users can select any page to revise—expertise plays no explicit role in such revisions. It has become the world’s largest “collective intelligence” experiment and one of the largest human projects ever to bring information into one source. The website receives enormous attention, with over eight billion page views per month in the English language, and over 500 million unique visitors per month.<sup>10</sup>

Contributions come from tens of millions of dedicated contributors who participate in an extensive set of formal and informal roles.<sup>11</sup> Some of these roles entail specific responsibilities in editing tasks; however, the Wikimedia Foundation employs a limited set of people and largely does not command its volunteers. Rather it helps develop a number of mechanisms to govern the co-production process by volunteers (Kane and Fichman 2009; Te’eni 2009; Zhang and Zhu 2011, Hill 2017). All these voluntary contributors are considered editors on Wikipedia. The organization relies on contributors to discover and fix passages that do not meet the site’s content tenets, but no central authority tells contributors how to allocate editorial time and attention.

The reliance on volunteers has many benefits but comes with many drawbacks. Among the latter, there is a long-standing concern that interested parties attempt to rewrite Wikipedia to serve their own parochial interests and views. Despite the persistence of such concerns, there is little systematic evidence pointing in one direction or another. Available evidence on conflicts suggests that contributors who frequently work together do not get into as many conflicts as those who do not, nor do their conflicts last as long (Piskorski and Gorbatai 2013). Additional evidence suggests a taste for prosocial and reciprocal behavior among contributors also plays an important role in

---

<sup>10</sup> “Wikipedia vs. the small screen”. <http://www.nytimes.com/2014/02/10/technology/wikipedia-vs-the-small-screen.html? r=1>, assessed June 2016.

<sup>11</sup> See [https://en.wikipedia.org/wiki/Wikipedia:User\\_access\\_levels](https://en.wikipedia.org/wiki/Wikipedia:User_access_levels), accessed June 2016.

fostering long-lasting cooperation among them (Algan et al. 2013). While such behavior could lead to edits from contributors with different points of view, there is no direct evidence that it leads to more content that finds compromises between opposite viewpoints.

While the Wikipedia community tries to attract a large and diverse community of contributors, there is general recognition that it invites many slanted and biased views. Moreover, the openness of Wikipedia’s production model (e.g., allowing anonymous contributions) is subject to sophisticated manipulations of content by interested parties. So there is widespread acceptance of the need for constant vigilance and review.

A key aspiration for all Wikipedia articles is a “neutral point of view” or NPOV (e.g., Majchrzak 2009, Hill 2017). To achieve this goal, “conflicting opinions are presented next to one another, with all significant points of view represented” (Greenstein and Zhu 2012). In practice, when multiple contributors make inconsistent contributions, other contributors devote considerable time and energy debating whether the article’s text portrays a topic from a NPOV. Because Wikipedia articles face virtually no limits to their number or size<sup>12</sup>—due to the absence of any significant storage costs or any binding material expense, conflicts can be addressed by adding more points of view to articles, rather than by eliminating them (e.g., Stvila et al. 2008). Like all matters at Wikipedia, contributors have discretion to settle disputes on their own—no command comes from the center of the organization. The center offers a set of norms for the dispute resolution processes, which today can be quite elaborate, including the three-revert edit war rule, as well as rules for the intervention of arbitration committees and mediation committees. Administrators can also decide to freeze an article under contention.

### **3. Data and Summary Statistics**

A number of statistical challenges arise when measuring micro-behavior of segregated conversations. First, because both contributors and articles may be slanted and biased, we must take both into account when developing a yardstick to compare the contributor to the contribution. That yardstick must enable a quantifiable method for studying whether contributors select content

---

<sup>12</sup> Over time a de facto norm has developed that tends to keep most articles under six to eight thousand words. This arises as editorial teams debate and discuss the length of the article necessary to address the topic of the page. Of course, some articles grow to enormous lengths, and editor contributors tend to reduce their length by splitting them into sub-topics. Prior work (Greenstein and Zhu 2016) finds that the average Wikipedia article is shorter than this norm (just over 4,000 words), but the sample does include a few longer articles (the longest is over 20,000 words).

with a slant similar to their own slant. Second, the slant and bias of articles changes because contributors revise articles.<sup>13</sup> Thus, we need a method that measures the changes as the content of articles change. Third, contributors themselves may also change as they gain experience by editing more articles with slants and biases similar or different from their own. Hence, we need a way to measure the evolution of contributors, as well as of their contributions.

Following an approach pioneered in Greenstein and Zhu (2016), we develop a sample of articles from Wikipedia. We focus on broad and inclusive definitions of U.S. political topics, including all Wikipedia articles that include the keywords “Republican” or “Democrat.” We start by gathering a list of 111,216 relevant entries from the online edition of Wikipedia on January 16, 2011. Eliminating the irrelevant articles and those concerning events in countries other than the United States<sup>14</sup> reduces our sample to 70,305. Our sample covers topics with many debates over contestable knowledge, ranging from the controversial topics of abortion, gun control, foreign policy, and taxation, to the less disputed ones relating to minor historical and political events and biographies of regional politicians. We next collect the revision history data from Wikipedia on January 16, 2011, which yields 2,891,877 unique contributors.

To mitigate concerns about manipulating statistical procedures, we rely on a modification of an existing method, developed by Gentzkow and Shapiro (2010), for measuring slant and bias in newspapers’ political editorials.<sup>15</sup> For example, Gentzkow and Shapiro (2010) find that Democratic representatives are more likely to use phrases such as “war in Iraq,” “civil rights,” and “trade deficit,” while Republican representatives are more likely to use phrases such as “economic growth,” “illegal immigration,” and “border security.”<sup>16</sup> Similarly, we compute an index for the slant of each article from each source, tracking whether articles employ these words or phrases that appears to slant toward either Democrats or Republicans.

---

<sup>13</sup> This is a property that Greenstein and Zhu (2012) confirmed in their study of Wikipedia articles.

<sup>14</sup> The words “Democrat” and “Republican” do not appear exclusively in entries about U.S. politics. If a country name shows up in the title or category names, we then check whether the phrase “United States” or “America” shows up in the title or category names. If yes, we keep this article. Otherwise, we search the text for “United States” or “America.” We retain articles in which these phrases show up more than three times. This process allows us to keep articles on issues such as “Iraq War,” but drop articles related to political parties in non-U.S. countries.

<sup>15</sup> Gentzkow and Shapiro (2010) characterize how newspapers also use such phrases to speak to constituents who lean toward one political approach over another.

<sup>16</sup> Several studies have applied their approach in analyzing political biases in online and offline content (e.g., Greenstein and Zhu 2012; Jelveh et. al. 2014). In addition, although Budak et al. (2014) use alternative approaches to measure ideological positions of news outlets, their results are consistent with Gentzkow and Shapiro (2010).



Like Gentzkow and Shapiro (2010), we investigate whether Wikipedia articles use words or phrases favored more by Republican or Democratic members of Congress. Gentzkow and Shapiro (2010) select such phrases based on the number of times they appear in the text of the 2005 *Congressional Record*, and apply statistical methods to identify those phrases that separate Democrat and Republican representatives. Their approach rests on the notion that each group uses a distinct “coded” language to speak to its respective constituents.<sup>17</sup> Each phrase is associated with a cardinal value that represents the degree to which each word or phrase is slanted. After offering considerable supporting evidence, Gentzkow and Shapiro (2010) estimate the relationship between the use of each phrase and the ideology of newspapers, using 1,000 words and phrases to identify whether those newspapers’ views tend to be more aligned with Democrat or Republican ideologies. As shorthand we refer to these 1000 words and phrases as “code phrases.”

This approach has several key strengths in that it has passed many internal validity tests, avoids many subjective elements, and provides a general yardstick for measuring the bias of newspaper articles. The approach also is effective when examining political bias in articles in economic journals (Jelveh et al. 2014), which we believe can be transferred to the context of Internet articles. Wikipedia’s contributors are unlikely to have used this yardstick to target these words for editing, though they might have included or excluded them when endeavoring to represent or exclude a specific point of view. The method also leads to a quantifiable measure of “neutral,” because the numbers are additive for finding the total slant of an article, and the range of slants can be normalized at the mean. An article is deemed unslanted or unbiased either when it includes no code phrases from many opposing points of view or when its use of Republican and Democrat code phrases equal the same cardinal value.<sup>18</sup>

In general, just as there is no definitive way to measure the “true bias” of a newspaper article in Gentzkow and Shapiro (2010), there is no definitive way to measure the true bias of an online encyclopedia article. Our normalization is valid under the assumption that the underlying differences among the population of contributors do not change over the sample period, and the variance of observed slant around this mean is random. As we illustrate below, because the analysis

---

<sup>17</sup> See Table I in Gentzkow and Shapiro (2010) for more examples.

<sup>18</sup> Greenstein and Zhu (2016) find no evidence that these two types of unslanted articles differ in their underlying traits. Hence, in this paper we treat them as identical.

focuses on the pairing of the slant of contributor/contribution, the inferences will be robust to small changes in the normalization.

### 3.1. Measures

#### 3.1.1. Dependent variables

Contributor Slant. Every article on Wikipedia has a revision history that, for every edit, records a pre-edit and post-edit version. We compute the slant index for both the pre- and post-edit article versions, take the difference between the two, and use this difference in slant as the *slant change* resulting from this edit. In this way, we obtain the slant change of every edit. For sequential edits from the same contributor that happened consecutively and without anyone else editing between them, we treat the sequence of edits as one single edit in all our analysis. These consecutive edits tend to be highly correlated, or could be several parts of a complete contribution, such as where the contributors saved their work several times.

Next, we focus on individual contributors as the unit of analysis. For our research purposes, we need to identify the bias and slant of contributors on the basis of their online political ideologies. To do so, we identify and measure the types of changes they make to Wikipedia articles. For every edit in our data, we take the difference between the pre-edit and post-edit versions of the article to determine the slant change of this edit. We assign each edit to each contributor, and assign a slant value for each edit. Under the assumption that every contributor has one fixed type of slant, we compute the *Contributor Slant* as the average value of the slant index of this contributor.

A zero value of *Contributor Slant* means the user's edits either contain a balanced set of Republican/Democratic words (weighted by their cardinal values) or do not include any of the slanted phrases. A negative or positive value of *Contributor Slant* means the contributor is Democrat-leaning or Republican-leaning, respectively. In our sample, 2,678,626 out of 2,891,877 unique contributors (92.6%) have a zero contributor slant, and over 225 thousand contributors make at least one slanted contribution.

Contributor Slant by Year. In our first analysis we will assume contributors have the same slant over their lifetime, and in the second analysis we relax the constraint that contributors maintain the same type of slant over time. In the latter, we divide contributors' edits by year and for each

year use the same calculation as for *Contributor Slant*, that is, we compute the average slant change of all the edits a contributor has made within that year. If a contributor's numeric value for slant remains unchanged throughout the years, then his or her *Contributor Slant by Year* equals *Contributor Slant*.

*Contributor Category* and *Contributor Category by Year*. We create two categorical variables. Based on *Contributor Slant* we create *Contributor Category*, which takes the value of -1, 0, or 1, representing contributors with a slant two standard deviations below mean, in between, and above mean, respectively. *Contributor Category by Year* is the yearly version of *Contributor Category*.

### 3.1.2. *Explanatory Variables*

*Prior Article Slant* and *Prior Article Category*. *Prior Article Slant* denotes an article's slant before a particular edit. This variable is used as the explanatory variable to analyze the article's relationship with the next contributor's slant. We also create a categorical variable, *Prior Article Category*, by categorizing *Prior Article Slant* into -1, 0, and 1 for articles with slant two standard deviations below mean, in between, and above mean, respectively.

*Contributor Years*. For every edit in our sample, this is the number of years the contributor has been on Wikipedia before he or she made this edit. This time variable is used to analyze whether a contributor's slant changes over time.

### 3.1.3. *Moderating Variables*

*Average Bias of Articles Edited*. Numerically, an article's bias equals the absolute value of its slant. *Average Bias of Articles Edited* is the average bias of all the articles that a contributor has edited. This variable helps measure the contributor's online experiences and helps us identify the role of content bias on a contributor's slant change over time.

*Fraction of Extreme Articles Edited*. We use this variable to characterize the contents of the articles that contributors interact with during their online experiences. An article is defined as *extreme* if its slant is more than two standard deviations away from the mean. *Fraction of Extreme Articles*

*Edited* equals the ratio between the number of extreme articles that the contributor has edited and the total number of articles the contributor edited. Like *Average Bias of Articles Edited*, the variable, *Fraction of Extreme Articles Edited*, helps identify the role of content bias on contributors' slant change over time.

### 3.1.4. Control Variables

*Prior Article Length and Prior Refs.* Apart from the article slant, there are some other time-varying article-specific characteristics that may affect the selection of the type of contribution. For instance, articles that are longer may incorporate more viewpoints, which then, in turn, tends to attract more contributors. Also, Wikipedia requires citations from major third-party sources as references for its article content (often listed at the bottom of the page), so articles with more references are also more likely to incorporate more outside arguments or controversial views at the time. Articles with these characteristics may tend to attract certain types of contributors. To control for these influences, we measure the length of the articles using the number of words in an article prior to a certain edit, denoted by *Prior Article Length*, and we measure the number of the article's external references, denoted by *Prior Refs.* These variables are included in the regressions on the relationship between contributor slant and the prior article slant of the article that the contributor chooses to edit.

*Number of Edits.* As with articles, there are time-varying characteristics of contributors that may affect their slant change over time. One of them is the total number of edits that a contributor has made so far, since people who make more edits may be affected more by the online contents. We use *Number of Edits*, the total number of edits *to date* that the contributor has made on Wikipedia, to control for such influence when analyzing the effect of time on contributor slant changes.

## 3.2. Summary Statistics

Table 1 presents the distribution of types of contributors over ten years. When computing the number of Democratic, Republican, and Neutral contributors to Wikipedia each year, we count each user ID only once—even if the user contributes many times in a year. There are 2,891,877 unique contributors in our sample. As noted above, 92.6% have zero contributor slant. We define

a contributor as *active* if his or her total number of edits is distributed in the top 10% of all contributors' total number of edits, which in this case equals a total of no less than three contributions in our sample. Active contributors comprise 10% of contributors, but they make 74% of the contributions in the entire sample. In other words, most of the edits in the sample come from experienced contributors – these are the contributors who we expect to be savvy about reading the existing slant of the articles and responding to that slant. Furthermore, while the number of neutral contributors who contribute each year is more than ten times that of contributors who have a slant, the proportion of active contributors in the neutral slant group (15.9%) is much smaller compared to the proportion of active contributors in the other two groups (63.8% and 65.5%). In summary, slanted contributors are more active than neutral contributors, and much of the slanted content comes from contributors making many edits.

In Table 2, we provide summary statistics of all variables used in our analysis. The unit of analysis in this table is contributor-edits, and the total number of observations is 10,948,696. Edits from all contributors who have ever contributed to the articles in our sample are included in this table. While in Table 1 we summarize on the level of *contributors*, in Table 2 we focus on all the *edits* made by the contributors within the entire time period. The two tables together help develop a broad understanding of both who contributes and what they contribute to the articles.

In general, the average *Contributor Slant* in our sample is negatively close to zero, while the average *Contributor Category* is positively close to zero. The summary statistics indicate that (1) Democrat-leaning contributors are, on average, more slanted than Republican-leaning contributors, and (2) all article versions in our sample exhibit a Democrat-leaning slant, with similar absolute values of extreme slant on both ends. There is also substantial variation across article versions for each of the three control variable measures, and we use the logarithm of these three control variables in our models since they are highly skewed.

We summarize the distribution of contributors' total number of edits over the ten years using Figure 1. Our sample reflects the well-known skewness of contributions to Wikipedia. More than 75% of the contributors in our sample contributed only once in the entire ten-year period. 97.5% of the contributors contributed fewer than 10 times, averaging to less than one contribution per year. Only 1% of the contributors contributed more than 30 times in our sample.

## 4. Empirical Results

### 4.1. Contributors' Participation Pattern on Wikipedia

For every edit in our sample, we look at the relationship between the contributor's slant and the article's slant that he or she chooses to edit by using the following regression model:

$$\text{Contributor Slant}_j = \alpha_0 + \alpha_1 \text{Prior Article Slant}_{it} + X_{it}B + \sigma_i + \eta_t + \varepsilon_{it} . \quad (1)$$

The coefficient  $\alpha_1$  identifies whether the average contribution follows BOF or OA. Here,  $X_{it}$  is a vector of the article's characteristics and control variables,  $\sigma_i$  is an article fix effect to control for any fixed differences among articles (despite many potential changes over many years), and  $\eta_t$  is a year fixed effect to control for any common trend in media/macroeconomic shocks that may differentially affect articles of different years. As an alternative approach, we use *Contributor Category* as the dependent variable, with *Prior Article Category* as the explanatory variable.

In Table 3, we report estimation results of Equation (1) using Ordinary Least Square (OLS) regressions. For the sake of analyzing participant behaviors, we drop the first version of all articles in our sample, since we do not have a prior article slant and cannot observe OA or BOF effect for such contributions. This reduces the number of observations in the sample to 10,878,391 and the number of articles to 66,389. Unless pointed out otherwise, all analysis samples used later in this paper is the same as this sample.

Models (1) through (3) use *Contributor Slant* as the dependent variable. Model (1) includes only *Prior Article Slant* as the explanatory variable. Model (2) adds in control variables *Log (Prior Article Length)* and *Log(Prior Refs)*. Model (3) replicates Equation 1, with article- and year- fixed effects included. The coefficients on *Prior Article Slant* is negative and significant in all three models. This indicates that an increase in the article's slant is associated with a decrease in the slant of its next contributor; namely, when the article is more Republican-leaning, it tends to attract a more Democrat-leaning user as its next contributor. That is consistent with OA behavior.

Models (4)-(6) repeat the analyses in Models (1)-(3) but replace *Contributor Slant* with *Contributor Category* as the dependent variable, and replace *Prior Article Slant* with *Prior Article Category* as the explanatory variable. Again, we find that the coefficients for the categorical explanatory variable *Prior Article Category* is negative and significant in all cases, suggesting that

the slant category of the next contributor is significantly negatively correlated with the slant category of the prior article. Results are similar across models and in line with our findings from Models (1)-(3).

We also partition the contributors by their frequency of edits and examine whether core and peripheral contributors behave similarly in our sample. *Core* contributors are the *active* contributors in Table 1; i.e. the top 10% contributors in terms of each contributor’s total number of edits. *Peripheral* contributors are contributors who made only one edit in our sample, here represents 75.5% of all contributors.

Table 4 reports the regression results of Equation 1 based on subsamples of core contributors and peripheral contributors. Again, both types of contributors demonstrate a similar OA pattern in their participation behavior, with peripheral contributors showing greater magnitude of the effect compared to core contributors. The results still hold after controlling for year and article fixed effects in the regressions.<sup>19</sup>

To further illustrate the OA pattern in contributors’ online participation, we use multinomial logistic regressions on the relationship between *Contributor Category* and *Prior Article Category*, with control variables and fixed effects similar to the specifications in Equation 1.

In Table 5, we present the estimation results. Again, Model (1) includes only *Prior Article Category* as the explanatory variable. Model (2) adds in control variables *Log (Prior Article Length)* and *Log(Prior Refs)*. Model (3) includes fixed effects. We can see that the coefficients for *Prior Article Category* are all statistically significant and have opposite signs with the categorical dependent variable. Take the coefficients of *Prior Article Category* in Model (1) as an example. The coefficient for *Prior Article Slant* is 2.10 when the *Contributor Category* is -1, which leads to a 4.0% increase<sup>20</sup> in the probability of attracting a next contributor whose *Contributor Category* equals -1 when the article’s prior slant increases by 1. Compared to the baseline coefficient, this result shows that when a prior article’s slant moves to a Republican-leaning slant by one category, it is eight times more likely that it will attract a Democrat-leaning user as its next contributor. Similarly, the coefficients in Model (2) and (3) suggest that the increase in the probability of

---

19 Besides core and peripheral contributors, there is also a middle group that includes 14.5% of contributors in our sample. Contributors in this middle group demonstrate a similar OA pattern as contributors in the other two groups, with a magnitude of the OA effect inbetween that of the core and the peripheral contributors.

<sup>20</sup>  $\frac{e^{-5.11+2.07}}{1+e^{-5.11+2.07}+e^{-5.25-2.41}} - \frac{e^{-5.11}}{1+e^{-5.11+2.07}} = 0.0456 - 0.0058 = 0.0398.$

attracting a subsequent contributor with an opposite slant is even higher than it was without control variables or year fixed effects. Overall, the results continue to support our previous findings of a greater OA effect than BOF effect in contributors' online participation.

#### 4.2. *Do Contributions from Contributors Change Over Time?*

In the previous analysis, we have assumed that every contributor's slant is constant over time. We now relax that assumption, and examine how a contributor's slant changes over time. We estimate the following equation:

$$\text{Contributor Slant by Year}_{jt} = \beta_0 + \beta_1 \text{Contributor Years}_{jt} + Z_{it}B + \mu_j + \epsilon_{it} . \quad (2)$$

The coefficient  $\beta_1$  can help identify whether and how contributor slant changes over time. Here  $Z_{it}$  includes a contributor's characteristics and controls for time-varying differences among contributors, such as *Number of Edits*.  $\mu_j$  is a contributor fix effect. Because it is not possible to estimate  $\mu_j$ , a contributor fix effect, for contributors who make one contribution, the number of observations that enter the regression with contributor fix effect becomes smaller. We try estimates with and without this effect.

In Table 6, Models (1) through (4) use the absolute value of *Contributor Slant by Year* as the dependent variable. We take the absolute value to capture how far away the contributor slant is from neutral, regardless of its sign. Model (2) includes contributor fixed effect, and Model (4) includes both contributor fixed effect and contributor characteristics as control variables.

The estimated coefficients of *Contributor Years* in all models are negative and statistically significant. The result means that, overall the average Wikipedia contributor slant declines over time. The average contributor slant moves closer to neutral by 0.0002 for every additional year the contributor stays in the community.

Although we observe an overall decline in the bias of contributors over time (e.g., the year 2008 is a notable exception to the trend), one might argue that such a decline arises as an artifact of the dictionary of code phrases we use. We compute the slant measure in 2005, which may become less relevant over time. If this is the case, we would expect to see the contributor slant decline only after 2005. To test this, we exclude all the observations after 2005 from our sample



and re-run the above OLS regression to see how the absolute value of *Contributor Slant by Year* changes during these years. Again, the results show a significant negative relationship between contributors' slant and contributor years, indicating that the decline in contributor slant is not due to decreasing relevance of our slant measure.

In addition to looking at how the average contributor slant changes, we use Markov matrix to illustrate how slant composition of contributors evolves over time. This matrix, reported in Figure 2, is constructed as follows: First, we divide in half every contributor's time that he or she has been on Wikipedia. Then, we divide the direction of this contributor's edits by attaching values (-1, 0, 1) to negative slant, zero slant, and positive slant edits. Based on the sum of these values for the first half and the second half of this contributor's activity, we can categorize the contributor as Democrat, Neutral, or Republican: If the sum of all edits in one half is negative (positive), the contributor is a Democrat (Republican), respectively. And, if the sum of all edits in this half is zero, the contributor is Neutral. We do this for each half of every contributor's activity on Wikipedia and accumulate them to get the overall transition probabilities in the entire community. We find that, for both democratic-leaning and republican-leaning contributors in the first half, there is more than a 70% chance that they will move to Neutral in the second half of their activities. As a result, the community in general has a tendency of moving towards neutral.

Since it is more likely that contributor slant declines over time instead of remaining constant throughout the years, we next examine whether our findings of OA in contributor participation is still valid under the different contributor slant assumption. We repeat the OLS regressions utilized above by using *Contributor Slant by Year* as the explanatory variable. From the results in Table 7, we can see that, just as in Table 4, the coefficients for *Prior Article Slant* and *Prior Article Category* remain negative and statistically significant ( $p < 0.001$ ). Moreover, compared to those under the constant contributor slant assumption, the magnitudes of the estimated coefficients are actually larger when using *Contributor Slant by Year* as the dependent variable. The results provide further support for our previous findings that there exists a significant OA pattern in contributors' participation in Wikipedia.

#### 4.3. Do Contributors Learn From Their Editing Experiences?

We next investigate how a contributor's prior editing experiences affects the slant of his or her contribution. Equation (3) adds the average bias of prior edited articles for each contributor, *Average Bias of Articles Edited*, and interacts it with *Contributor Years*, yielding:

$$\text{Contributor Slant by Year}_{jt} = \gamma_0 + \gamma_1 \text{Contributor Years}_{jt} + \gamma_2 \text{Average Bias of Articles Edited}_{jt} + \gamma_3 \text{Contributor Years}_{jt} \times \text{Average Bias of Articles Edited}_{jt} + Z_{it}B + \mu_j + \tau_{it} . \quad (3)$$

The coefficient  $\gamma_3$  estimates the moderating effect of extreme contents on contributors' slant change over time. Like Equation (2),  $Z_{it}$  refers to *Number of Edits*, which is a contributor characteristics variable controlling for time-varying differences among contributors, and  $\mu_j$  is a contributor fix effect to control for any fixed differences among contributors. In an alternative specification we also use *Fraction of Extreme Articles Edited* as an alternative measure for extreme contents, including this variable and its interaction term with *Contributor Years* in Equation 3.

Regression results using each of the two content measures are reported in parallel in Table 8. Model (1) and Model (2) estimate the moderating effect of *Average Bias of Articles Edited*. The coefficients for the interaction terms are negative and statistically significant, which indicates that if a contributor has been interacting with articles that are very biased, his or her own slant becomes neutral more quickly over time. The estimated coefficients show that the average article bias does have a significant influence on contributors' slant change. Models (3)-(4) replaces *Average Bias of Articles Edited* with *Fraction of Extreme Articles Edited*. Again, the estimated coefficients of the interaction terms are negative and statistically significant. However, the findings also are mildly mixed because the coefficients for *Contributor Years* are near zero, and change sign with different specifications.

#### 4.4. Rate of Slant Change: How Long Will It Take for Contributors to Become Neutral?

The presence of considerable heterogeneity makes it challenging to characterize the implications of the patterns of these findings. Having observed the tendency of contributor slant change over time, we next estimate how long it takes for a contributor's slant to gradually converge to neutral if this tendency continues.

We use a Markov Chain Process to simulate the slant convergence. Although a contributor's slant exhibits long-term trend over the years, it fluctuates frequently, and this should be accounted for. We divide slant into different bins and investigate how a contributor's slant changes from one bin to another. *Contributor Slant by Year* is divided into seven bins, divided by the  $\pm 0.5$ ,  $\pm 1.5$ , and  $\pm 2.5$  standard deviations intervals. The middle bin represents a neutral slant; the first and last bins represent extreme slants. We then compute a transition matrix for contributor slant based on our empirical data: For each year, we compute the proportions of contributors whose yearly slant moves from one slant bin to another, and fill the probabilities in the transition matrix for this year. Averaging the transition matrices among all years gives us the final transition matrix we use in our simulation, reported in Figure 3.

In this transition matrix, the rows denote the starting bins and the columns denote the ending slant. Bin 4 represents a neutral slant, defined as a slant index ranging from -0.5 to 0.5 standard deviations away from the mean. We find that: (1) the probabilities on the diagonal are relatively large. As expected, contributors tend to have a higher chance of staying near their original slant; and (2) the farther the end bins are from the start bins, the smaller the probabilities. This indicates that contributor slant change is a gradual and accumulative process, and it is not likely that the contributor's slant would suddenly jump from one extreme to another.

Next, we use the transition matrix to simulate the contributor slant change process over time (see Table 9). We compute the time it takes for a contributor to have a greater than 50% probability of moving to neutral. As expected, the length of time depends on the contributor's original slant: Extremely slanted contributors spend a longer time moving to neutral than slightly slanted contributors. More surprisingly, we find that on average, it takes one more year for the Republicans to become neutral than for Democrats.

We test for several possible reasons why Republican contributors converge to neutral slant slower than Democratic contributors. First, it could be that Republican contributors in general display more BOF behavior than Democratic contributors. Regression results of Equation (1) using the two groups respectively do not support this explanation. In fact, Republican contributors in general show stronger magnitude of OA compared to Democratic contributors.

Second, Republican contributors might choose to edit less extreme articles compared to Democratic contributors, so that they are less influenced during their interaction with online content. However, we find no statistically significant difference between the level of content

extremeness for the articles edited by Republicans or Democrats. The distributions contain similar bias and variance.

A third possible reason might stem from the contributors' numbers of edits – that is, Republican contributors make fewer edits in our sample than Democrats, so their experience has less of an effect on the overall tendency, and may differ in some way. Summary statistics provide evidence for this explanation. In our sample, the total number of edits from Democratic contributors is about 1.5 times that from Republican contributors.

Furthermore, the two types of contributors examine different topics, and each of these display different OA/BOF behavior. We characterize the heterogeneity of OA/BOF among different topics, using Wikipedia's classification for articles. We create dummy variables for each topic categories and modify Equation (1), adding these dummies and their interactions with *Prior Article Slant*. We then compute the OA effect for each topic category using the regression results. There are 24 categories of topics in the sample, and these are not mutually exclusive. Articles can speak to one or more than one topic, and these rarely change over the lifetime of an article. We estimate this modification to Equation (1) for the entire sample, and for two mutually exclusive sub-samples, one consisting of Republican contributors and one for Democrat Contributors. We report the results in Table 10.

Consistent with our overall findings, the majority of topics display OA for contributors from both parties. For example, the four topics with the most edits – Foreign Policy, Government, War and Peace, and Biographies – display an overall pattern of OA. Most interesting are the departures from this pattern. Among the ten topics receiving the most edits, three topics – Budget and Economy, Civil Rights, and Crime – display OA overall, with either Democrats displaying BOF and Republicans displaying OA, or no significant pattern. This can happen if the Democratic contributors resist changing content when Republicans try to insert their point of view. The same pattern in the opposite direction, with Democrats displaying OA and Republicans displaying BOF, occurs only on one topic with much fewer edits—Healthcare. Three topics—Homeland Security, Energy, and Tax—display evidence of a segregated conversation, where both parties engage in BOF, and they are not in the top ten in terms of the number of edits. In these three topics, however, the BOF effect of Republican contributors is much stronger than that of Democrats, indicating that Republicans' edits are the relatively stronger force that contributes to these segregated conversations.

Overall, Table 10 suggests Republican and Democratic contributors do have different experiences, selecting among different groups of articles to edit, most frequently those with a different viewpoint. The weight of experience results in OA overall, with Republican editors experiencing (somewhat) segregated conversations less frequently (as a numerical matter). To say it another way, Republicans converge more slowly to neutral because of the proportion of time they find themselves on the opposite side of the content—in comparison to Democrats. In sum, the findings again support our primary conclusions that (1) online experiences change contributors’ slant and (2) there is a tendency for Wikipedia contributors’ slants to converge.

## 5. Robustness of Findings and Alternative Explanations

We further corroborate our findings by performing the following robustness tests.

### 5.1. *Is the Measure of Contributor Slant Representative of Ideologies?*

First, since the measure of contributors’ political ideologies and slant are computed entirely on the basis of data from Wikipedia, one might be concerned about whether such a slant measure is representative of contributors’ real-world political ideologies. Also, a neutral article in our sample can either be interpreted as having no slanted words at all or as having equal numbers of very slanted words. These concerns might lead to questioning the external validity of the slant measure.

To address this concern we use an alternative measure of slant and bias of contributors. We match the voting data from the 2000 Presidential Election to locations affiliated with IP addresses of contributors.<sup>21</sup> Because Wikipedia only reveals IP addresses for contributors without user IDs, we restrict our sample to contributors who are not logged in when editing the articles and also drop contributors whose IP addresses indicate that they are located outside the United States. Using OLS regressions, we then test the relationship between the voting record and *Contributor Slant*. Note that this analyzes the behavior of a different population of contributors than the contributors we have examined thus far.<sup>22</sup> This regression is valid under the assumption that a contributor has – on average – the political tastes of the regions from which they live.

---

<sup>21</sup> The data on geolocation of IP comes from MaxMind. We match on county records.

<sup>22</sup> The identifies of contributors are known after they register, and when they edit after logging on. An anonymous edit comes from either an unregistered contributor or from an editor who choses not to logon before editing. Hence, it is possible for the samples to include some of the same contributors, but it is not possible to know what fraction.

Table 11 presents the results. *RepPerc* denotes the percentage of Republican votes in the contributor's county. As we use positive values in the slant index to indicate Republican-leaning ideologies for Wikipedia users and articles, the positive and statistically significant coefficient of *RepPerc* suggests that a user's *Contributor Slant* index is larger when the county from which he or she votes has a higher percentage of Republican votes. The results are qualitatively similar to the prior estimates. This provides evidence that the measure of contributors' slant reflects contributors' real world political ideologies.

## 5.2. *What Else Could Be Driving the OA behavior?*

The effect of OA in contributors' voluntary editing behavior indicates that contributors are more likely to edit articles with the opposite slant. However, apart from the interpretation of contributors being attracted by the article slant, this could also be due to a "correcting" behavior between contributors, which might have little to do with the article's slant. On Wikipedia, we sometimes see edits that are reverted and added back within a short time, which are called "edit wars." Could these edit wars be driving the OA effect? We address this question by including only the initial edits of every contributor when they revise an article for the first time. Doing so rules out edit wars or any possible correcting behavior later in the edits.

We observe from Table 12 that the signs and statistical significance of the estimated coefficients do not change, and the magnitude of the coefficients becomes even larger, indicating an even stronger OA effect than when investigating all the edits. The results further strengthen the robustness of the OA effect.

We also conduct several additional robustness checks to make sure the OA effect is not driven by alternative explanations. First, our slant index is measured on the basis of frequently used phrases, or code phrases, favored by party representatives. It may be the case that longer articles tend to contain more code phrases and are therefore more measurable. In this case, long articles could drive our results. To rule out this explanation, we eliminate outlying long articles from our full sample, that is, articles that are more than two standard deviations above the mean article length. We obtain similar results.

Second, since we measure article slant using code phrases, the articles whose titles contain code phrases might tend to show greater biases in our sample simply because these code phrases are more likely to be used repetitively in the article content. To check the robustness of our finding,

we exclude from our sample all articles whose title contain code phrases, which is 1.77% of all articles. Again, we find a significant OA effect from the results.

Third, it is possible that certain code phrases are chosen simply because these words do not have other commonly-used synonyms that are neutral or of the opposite slant. In this case, as our measure captures the contributor's choice of words describing the same concept for a given topic, one's contribution may be slanted merely because he or she could not find neutral substitutes of the code phrases to choose from. We rely on the experiences of a legal and copyediting professional to identify these instances in our dictionary and leave only code phrases with natural substitutes. After re-measuring the slant index for articles and contributors, we repeat our analyses and find no significant change in our results. Therefore, the OA effect is not driven by instances where contributors do not have a choice for substitute phrases.

Finally, we test if the OA effect is driven only by extremely slanted articles, or if the finding is universal among all articles. We eliminate from our full sample articles with slant index two standard deviation points away from the mean. Changing this threshold to articles without slant in the top and bottom 10% does not differ qualitatively in results. The estimated coefficients with subsamples have the same signs but larger absolute values.

### *5.3. Could There Be Vintage Effects Among Contributors?*

Perhaps the average contributor slant declines over years because of the differences among people joining Wikipedia in different years. That is, there may exist some pattern of user vintage effects across the years. For instance, compared to people who contributed later, those who contributed when Wikipedia was still in its early stage may not have been as proficient in editing neutral content as those who entered later. In this case, we may see that contributors who entered earlier are more slanted, and contributors who entered later are more neutral, on average.

We compute the average slant of contributors entering in different years and plot the results in Figure 4. As we can see, there is no obvious inclining or declining pattern in the average contributor slant across the years. Contributors who entered earlier are not systematically more neutral, nor are they more slanted, compared to those who entered later. This shows there are no vintage effects influencing the contributor slant convergence tendency in our findings. This finding also suggests that the change in slant over time is not caused by entry and exit of contributors exhibiting extreme bias.

## 6. Conclusion and Discussion

This research shows that Wikipedia has a remarkable record of bringing opposing opinions into the same conversation through examining two micro-behaviors of contributors, the target and evolution of their contribution. Our findings point toward patterns that lead contributors to offer content to those with different points of view, which we call the OA effect. We also show that contributors moderate their contribution over time. The change in contributions is especially large for contributors who interact with articles that are more extreme and have greater biases. These effects reinforce the prevalence of unsegregated conversations at Wikipedia over time. We also estimate that this slant convergence process takes one year longer on average for Republicans than for Democrats. In summary, we find that the majority of Wikipedia's contributors do not segregate into a conversation that excludes other viewpoints. Contributors interact with those of opposite viewpoints much more frequently than they silo themselves and participate in echo chambers.

Our findings have important implications for both theoretical research and practice. We offer a two-step method for identifying the mechanisms contributing to polarization that distinguishes selection from evolution. Nothing in these methods presumes the results; the method can flexibly measure contributions to (un)segregated conversations in a variety of settings.

These findings have implications for when collective intelligence is hampered by the enthusiasm or frenzy of a crowd. Collective intelligence should be more trustworthy when mechanisms encourage confrontation between distinct viewpoints. It also should adopt processes, as in Wikipedia, which retains contributors who learn to moderate their contributions from their experience.

It is not as if Wikipedia avoids its share of disagreements and confrontations, so the findings also raise a subtle question: How does Wikipedia transform controversial topics into arguments that include many points of view and sustain the community over time? We believe that this success arises from the institutions that help overcome the challenges affiliated with aggregating contested knowledge. For one, the aspiration of achieving NPOV directs attention to specific areas. No side can claim exclusive rights to determine the answer, which allows every contributor to add another paragraph if it diffuses an issue by giving voice to dissent. In addition, miniscule storage and transmission costs reduce the cost of listing another view on a web page. Our results also suggest that the conflict resolution mechanisms and the mix of informal and formal norms at



Wikipedia play an essential role in encouraging a community that works towards a neutral point of view. We believe future work can compare alternative norms and mechanisms and help inform design of information aggregation mechanisms in online platforms.

These findings also raise questions for the market design literature about how the structure of interaction between contributor and content on other online social media – such as Facebook, Twitter, and Reddit – shapes the prevalence of (un)segregated conversations. We speculate that some simple design differences may have profound consequences for (un)segregating conversations. For example, Wikipedia contributors can both add material and remove material or refine the content in myriad ways, whereas contributors on Facebook/Twitter only add additional content on top of what is already there. Allowing for removing or editing anyone’s contributions can change how the reader and writer choose to direct the conversations, resulting in contributions from different points of view. Some platforms also aggregate contributions in ways that shape the prevalence of segregation. For example, on Yelp (e.g., rating restaurants) or Rotten Tomatoes (e.g., rating movies) additional material can be added without limit, the platform provides a numerical summary that can direct conversations between readers and reviewers. Our results frame questions about whether a numerical summary motivates others with views that differ from the summary or attracts more reviews from those who agree with it.

These findings also highlight the importance of platform design in social media. For example, on Facebook, an algorithm selects content for users, and its design increases the chance that participants read and write contents only in a community of like-minded people. Segregated conversation is also more likely on Facebook or Twitter due to processes that reinforce birds of feature to stick together. After all, a user often only sees content from his or her friends. Wikipedia contributors have the option to be exposed to different opinions and can freely make the choice of reading and writing any content on the platform. Future work can focus on the heterogeneous effect of online participation on different contributor subgroups—for example, with interest in different political topics, or participation in different types of online platforms, such as resource-sharing platforms versus communities of innovation. In addition, existing literature on open communities investigates the content production more frequently than the contributors themselves. Given the huge number of volunteers on Wikipedia, as well as the enormous attention this community gets from around the globe, we hope to see more research on Wikipedia’s online participation and interactions, as well as on the mechanisms behind changes to its content.

## References

- Ahn, Y-Y., Han, S., Kwak, H., Moon S., and Jeong H. 2007. "Analysis of topological characteristics of huge online social networking services." In *Proceedings of the 16th international conference on World Wide Web (WWW'07)*. ACM, New York, NY, USA, 835–844.
- Algan, Y., Benkler, Y., Morell, M.F. and Hergueux, J. 2013, July. "Cooperation in a peer production economy: Experimental evidence from Wikipedia." *Workshop on Information Systems and Economics*.
- Arazy, O., Nov, O., Patterson, R., and Yeo, L. 2011. "Information quality in Wikipedia: The effects of group composition and task conflict." *Journal of Management Information Systems* 21(4): 71–98.
- Arazy, O., Ortega, F., Nov, O., Yeo, L., and Balila, A. 2015. "Functional roles and career paths in Wikipedia." *Computer Supported Cooperative Work (CSCW)*: 1092–1105.
- Barsade, S. G. 2002. "The ripple effect: Emotional contagion and its influence on group behavior." *Administrative Science Quarterly* 47(4): 644–675.
- Benkler, Y. 2006. *The Wealth of Networks: How Social Production Transforms Markets and Freedom*. Yale University Press.
- Bernhardt D., Krasa S., and Polborn, M. 2008. "Political polarization and the electoral effects of media bias." *Journal of Public Economics* 92(5-6): 1092–1104.
- Besley, T., and Prat, A. 2006. "Handcuffs for the grabbing hand? Media capture and government accountability." *American Economic Review* 96(3): 720–736.
- Budak, C., Goel, S., and Rao, J. M. 2014. "Fair and balanced? Quantifying media bias through crowdsourced content analysis." Working paper. SSRN: <http://ssrn.com/abstract=2526461>.
- Burke, M. and Kraut, R. 2008. "Mopping up: Modeling Wikipedia promotion decisions." *Proceedings of the 2008 ACM Conference on Computer Supported Cooperative Work (CSCW)*: 27–36.
- Carr, N. 2008. *The Big Switch: Rewiring the World, from Edison to Google*. Norton, New York.
- Chiang, C.F. and Knight, B. 2011. "Media bias and influence: Evidence from newspaper endorsements." *Review of Economic Studies* 78(3): 795–820.
- Chiu, C-M., Hsu, M-H., and Wang, E. 2006. "Understanding knowledge sharing in virtual communities: An integration of social capital and social cognitive theories." *Decision support systems* 42(3): 1872–1888.
- Collier, B., Burke, M., Kittur, N., and Kraut, R., 2008. "Retrospective versus prospective evidence for promotion: The case of Wikipedia." *2008 Meeting of the Academy of Management*.
- DellaVigna, S. and Kaplan, E. 2007. "The Fox News effect: Media bias and voting." *Quarterly Journal of Economics* 122(3): 1187–1234.
- Durante, R. and Knight, B. 2012. "Partisan control, media bias, and viewer responses: Evidence from Berlusconi's Italy." *Journal of the European Economic Association* 10(3): 451–481.
- Forte, A., Kittur, N., Larco, V., Zhu, H., Bruckman, A., and Kraut, R. E. 2012. "Coordination and beyond: Social functions of groups in open content production." *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work*: 417–426.
- Forte, A., Larco, V., and Bruckman, A. 2009. "Decentralization in Wikipedia governance." *Journal of Management Information Systems* 26(1): 49–72.

- Gentzkow, M., and Shapiro, J.M. 2003. "Media, education, and anti-Americanism in the Muslim world." *Journal of Economic Perspectives* 18(3): 117–133.
- Gentzkow, M., and Shapiro, J.M. 2010. "What drives media slant? Evidence from U.S. daily newspapers." *Econometrica* 78(1): 35–71.
- Gentzkow, M., and Shapiro, J.M. 2011. "Ideological segregation online and offline." *Quarterly Journal of Economics* 126(4): 1799–1839.
- Greenstein, S., and Zhu, F. 2012. "Is Wikipedia biased?" *American Economic Review Papers and Proceedings* 102(3): 343–348.
- Greenstein, S., and Zhu, F. 2016. "Open content, Linus' Law, and neutral point of view." *Information Systems Research* 27(3): 618–635.
- Groseclose, T. and Milyo, J. 2005. "A measure of media bias." *Quarterly Journal of Economics* 120(4): 1191–1237.
- Halfaker, A., Kittur, A., and Riedl, J. 2011. "Don't bite the newbies: How reverts affect the quantity and quality of Wikipedia work." *Proceedings of the 7th International Symposium on Wikis and Open Collaboration (WikiSym'11)*: 163–172.
- Hill, M. 2017. "Almost Wikipedia: Eight early Encyclopedia projects and the mechanisms of collective action." Working paper. <https://mako.cc/academic/hill-almost-wikipedia-DRAFT.pdf>, accessed January, 2017.
- Janis, I. L. 1982. "Groupthink: Psychological studies of policy decisions and fiascoes." *Houghton Mifflin*.
- Jelveh, Z., Kogut, B., and Naidu, S. 2014. "Political language in economics." *Columbia Business School Research*, Paper No. 14-57.
- Jemielniak, D. 2014. *Common Knowledge?: An Ethnography of Wikipedia*. Palo Alto, CA: Stanford University Press.
- Jeppesen, L.B., and Frederiksen, L. 2006. "Why do users contribute to firm-hosted user communities? The case of computer-controlled music instruments." *Organization Science* 17(1): 45–63.
- Kane, G. C., and Fichman, R. G. 2009. "The shoemaker's children: Using Wikis for information systems teaching, research, and publication." *MIS Quarterly* 33(1): 1–17.
- Kittur, A., Chi, E., Pendleton, B. A., Suh, B., and Mytkowicz, T. 2007. "Power of the few vs. wisdom of the crowd: Wikipedia and the rise of the bourgeoisie." *World Wide Web* 1(2): 19.
- Kittur, A., and Kraut, R. 2008. "Harnessing the wisdom of crowds in Wikipedia: Quality through coordination." *ACM Conference on Computer Supported Cooperative Work*: 37–46.
- Kittur, A., Suh, B., Pendleton, B. A., and Chi, E. H., 2007. "He says, she says: Conflict and coordination in Wikipedia." *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*: 453–462.
- Larcinese, V., Puglisi, R., and Snyder, J.M. 2007. "Partisan bias in economic news: Evidence on the agenda-setting behavior of U.S. newspapers." *Journal of Public Economics* 95(9): 1178–1189.
- Laster, D., Bennett, P., and Geoum, I.S. 1999. "Rational bias in macroeconomic forecasts." *Quarterly Journal of Economics* 114 (1): 293–318.
- Lawrence, E., Sides, J., and Farrell, H. 2010. "Self-segregation or deliberation?: Blog readership, participation, and polarization in American politics." *Perspectives on Politics* 8(1): 141–157.
- Lee, Y.-J., Hosanagar, K., and Tan, Y. 2015. "Do I follow my friends or the crowd? Information cascades in online movie ratings." *Management Science* 61(9): 2241–2258.

- Ma, M., and Agarwal, R. 2007. "Through a glass darkly: Information technology design, identity verification, and knowledge contribution in online communities." *Information Systems Research* 18(1): 42–67.
- Majchrzak, A. 2009. "Where is the theory in Wikis?" *MIS Quarterly* 33(1): 18–20.
- McPherson, M., Smith-Lovin, L. and Cook, J.M. 2001. "Birds of a feather: Homophily in social networks." *Annual Review of Sociology*: 415–444.
- Muchnik, L., Aral, S., and Taylor, S.J. 2013. "Social influence bias: A randomized experiment." *Science* 341(6146): 647–651.
- Mullainathan, S., and Shleifer, A. 2005. "The market for news." *American Economic Review* 95(4): 1031–1053.
- Page, S, 2007. *The Difference. How the Power of Diversity Creates Better Groups, Firms, Schools and Societies*. Princeton University Press; Princeton, NJ.
- Piskorski, M.J., and Gorbatai, A. 2013. "Testing Coleman's social-norm enforcement mechanism: Evidence from Wikipedia." Working paper, Harvard Business School, Boston.
- Purdy, W. 2015. "Radicalization: Social media and the rise of terrorism." *Testimony presented before the House Committee on Oversight and Government Reform's Subcommittee on National Security*, available at <https://oversight.house.gov/wp-content/uploads/2015/10/10-28-2015-Natl-Security-Subcommittee-Hearing-on-Radicalization-Purdy-TRC-Testimony.pdf>
- Reuter, J., and Zitzewitz, E. 2006. "Do ads influence editors? Advertising and bias in the financial media." *Quarterly Journal of Economics* 121(1): 197–227.
- Rosenblat, T.S. and Mobius, M.M. 2004. "Getting closer or drifting apart?" *Quarterly Journal of Economics* 119(3): 971–1009.
- Schroeder, A., and Wagner, C. 2012. "Governance of open content creation: a conceptualization and analysis of control and guiding mechanisms in the open content domain." *Journal of the American Society for Information Science and Technology* 63(10): 1947–1959.
- Stone, D.F. 2011. "Ideological media bias." *Journal of Economic Behavior & Organization* 78(3): 256–271.
- Stvilia, B., Twidale, M. B., Smith, L. C., and Gasser, L. 2008. "Information quality work organization in Wikipedia." *Journal of the American Society for Information Science and Technology* 59(6): 983–1001.
- Sunstein, C.R. 2001. *Echo Chambers: Bush v. Gore, Impeachment, and Beyond*. Princeton University Press.
- Te'eni, D. 2009. "Comment: The Wiki way in a hurry—The ICIS anecdote." *MIS Quarterly* 33(1): 20–22.
- Van Alstyne, M., and Brynjolfsson, E. 2005. "Global village or cyber-balkans? Modeling and measuring the integration of electronic communities." *Management Science* 51(6): 851–868.
- Xu, S.X., and Zhang, X.M. 2013. "Impact of Wikipedia on market information environment: Evidence on management disclosure and investor reaction." *MIS Quarterly* 37(4): 1043–1068.
- Zhang, X.M., and Zhu, F. 2011. "Group size and incentives to contribute: A natural experiment at Chinese Wikipedia." *American Economic Review* 101(4): 1601–1615.
- Zitzewitz, E. 2001. "Measuring herding and exaggeration by equity analysts and other opinion sellers." Working Paper, Stanford Graduate School of Business. <http://www.dartmouth.edu/~ericz/chapter1.pdf>.

TABLE 1:  
Distribution of different types of contributors over years

Year	# Democrat Contributors Contributed	# of Active Democrat Contributors	# Republican Contributors Contributed	# of Active Republican Contributors	# Neutral Contributors Contributed	# of Active Neutral Contributors
2001	211	145	160	100	429	79
2002	434	327	420	324	3,510	767
2003	1,277	970	1,318	1,031	12,356	2,742
2004	5,191	3,844	5,170	3,932	56,506	11,580
2005	17,009	11,341	16,274	11,127	208,838	37,760
2006	33,512	20,786	33,106	21,004	517,820	85,756
2007	37,178	22,673	36,870	23,125	632,147	97,213
2008	33,517	20,121	33,803	20,786	573,551	89,220
2009	24,907	16,233	24,963	16,812	476,385	74,391
2010	19,434	12,974	19,518	13,454	422,711	60,920
2011	2,561	2,298	2,914	2,660	21,411	5,220
Total	175,231	111,712	174,516	114,355	2,925,664	465,648

Notes:

Definition: “# Democrat/Republican/Neutral contributors contributed” is the total number of contributors with negative/zero/positive *Contributor Slant* that have ever contributed to the articles in our sample.

Definition: “# of active Democrat/Republican/Neutral contributors” is the number of “Democrat/Republican/Neutral contributors contributed” whose total number of edits is distributed in the top 10% of all contributors’ total number of edits.

Final year, 2011, is sampled in January, which accounts for the low numbers in that year.

TABLE 2:  
Summary statistics of variables used in the main analyses

Variable	Mean	Std. dev.	Min	Max
Contributor Slant	-0.00025	0.02451	-1.22873	0.99807
Contributor Category	0.00077	0.11400	-1	1
Prior Article Slant	-0.05678	0.20786	-0.60507	0.62365
Prior Article Category	-0.05726	0.26403	-1	1
Prior Article Length	4049.76	3851.61	0	1,963,441
Prior Refs	33.9830	60.9042	0	1,636
Contributor Slant by Year	-0.00003	0.02361	-1.22873	0.99807
Contributor Category by Year	0.00118	0.12051	-1	1
Contributor Years	1.04022	1.36555	0.00274	9.79726
Number of Edits	1175.72	7567.79	1	122,264
Average Bias of Articles Edited	0.13759	0.11257	0	0.62365
Fraction of Extreme Articles Edited	0.07461	0.17640	0	1
RepPerc	0.45236	0.14383	0.093	0.919

Notes: Number of observations in this table is 10,948,696.

TABLE 3:  
OLS Regressions on the Relationship between Contributor Slant and Prior Article Slant

Model	(1)	(2)	(3)	(4)	(5)	(6)
Dependent Variable	Contributor Slant	Contributor Slant	Contributor Slant	Contributor Category	Contributor Category	Contributor Category
Prior Article Slant	-0.0075*** [0.0001]	-0.0074*** [0.0001]	-0.0167*** [0.0004]			
Prior Article Category				-0.0123*** [0.0002]	-0.0124*** [0.0002]	-0.0197*** [0.0009]
Log(Prior Article Length)		0.0005*** [0.0000]	0.0009*** [0.0001]		0.0014*** [0.0000]	0.0017*** [0.0003]
Log(Prior Refs)		-0.0003*** [0.0000]	-0.0009*** [0.0001]		-0.0008*** [0.0000]	-0.0024*** [0.0004]
Observations	10,878,391	10,878,391	10,878,391	10,878,391	10,878,391	10,878,391
Adjusted R-squared	0.005	0.006	0.006	0.001	0.001	0.001
Year FE	No	No	Yes	No	No	Yes
Article FE	No	No	Yes	No	No	Yes
Number of Articles	66,389	66,389	66,389	66,389	66,389	66,389

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. Observations in this panel are all the edits of the Wikipedia articles in our sample from 2001 to 2011. *Contributor Slant* is defined as the average slant change of all edits a contributor has made on Wikipedia. *Prior Article Slant* is the slant of the article before a particular edit. *Log(Prior Article Length)* is the logarithm of the article's total number of words. *Log(Prior Refs)* is the logarithm of the number of external references in the article plus one.

TABLE 4:  
OLS Regressions on the Relationship between Contributor Slant and Prior Article Slant, Core vs. Peripheral Contributors

Sample	Core Contributors	Peripheral Contributors	Core Contributors	Peripheral Contributors	Core Contributors	Peripheral Contributors	Core Contributors	Peripheral Contributors
Model	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Dependent Variable	Contributor Slant	Contributor Slant	Contributor Slant	Contributor Slant	Contributor Category	Contributor Category	Contributor Category	Contributor Category
Prior Article Slant	-0.0021*** [0.0000]	-0.0211*** [0.0002]	-0.0056*** [0.0002]	-0.0497*** [0.0012]				
Prior Article Category					-0.0063*** [0.0002]	-0.0237*** [0.0004]	-0.0109*** [0.0006]	-0.0410*** [0.0014]
Log(Prior Article Length)			0.0005*** [0.0000]	0.0035*** [0.0004]			0.0009*** [0.0001]	0.0055*** [0.0017]
Log(Prior Refs)			-0.0004*** [0.0000]	-0.0025*** [0.0003]			-0.0016*** [0.0002]	-0.0043*** [0.0011]
Observations	8,019,333	2,180,327	8,019,333	2,180,327	8,019,333	2,180,327	8,019,333	2,180,327
Adjusted R-squared	0.001	0.014	0.003	0.016	0.000	0.002	0.001	0.002
Year FE	No	No	Yes	Yes	No	No	Yes	Yes
Article FE	No	No	Yes	Yes	No	No	Yes	Yes
Number of Articles	66,313	46,856	66,313	46,856	66,313	46,856	66,313	46,856

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Definition: “Core Contributors” is the same as “active” contributors in Table 1; i.e. contributors whose total number of edits is distributed in the top 10% of all contributors’ total number of edits.

Definition: “Peripheral contributors” are contributors who made only 1 edit in our sample.



TABLE 5:  
Logit Regressions on the Relationship between Contributor Category and Prior Article Category

Model	(1)		(2)		(3)	
Dependent Variable	Contributor Category=-1	Contributor Category=1	Contributor Category=-1	Contributor Category=1	Contributor Category=-1	Contributor Category=1
Prior Article Slant	2.0743*** [0.0266]	-2.4063*** [0.0135]	2.0819*** [0.0269]	-2.3404*** [0.0133]	2.1042*** [0.0270]	-2.2918*** [0.0132]
Log(Prior Article Length)			-0.0344 [0.0045]	0.1486*** [0.0052]	-0.0115 [0.0051]	0.1859*** [0.0058]
Log(Prior Refs)			-0.2232*** [0.0032]	-0.3128*** [0.0030]	-0.2851*** [0.0042]	-0.4079*** [0.0040]
Year FE	No		No		Yes	
Article FE	No		No		Yes	
Observations	10,878,391		10,878,391		10,878,391	
Pseudo R-squared	0.021		0.038		0.043	

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

TABLE 6:  
Regressions of Contributor Slant Change over the Years

Model	(1)	(2)	(3)	(4)
Dependent Variable	Abs(Contributor Slant by Year)	Abs(Contributor Slant by Year)	Abs(Contributor Slant by Year)	Abs(Contributor Slant by Year)
Contributor Years	-0.0009*** [0.0000]	-0.0002*** [0.0000]	-0.0002*** [0.0000]	-0.0002*** [0.0000]
Log(Number of Edits)			-0.0005*** [0.0000]	-0.0001*** [0.0000]
Observations	10,878,391	10,878,391	10,878,391	10,878,391
R-squared	0.003	0.003	0.004	0.004
Contributor FE	No	Yes	No	Yes

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. Observations in this panel are the edits made by contributors. The dependent variable *Contributor Slant by Year* denotes the contributor's slant measured on the basis of the edits made within that year. *Contributor Years* denotes the number of years the contributor has been on Wikipedia. *Log(Number of Edits)* is the logarithm of the amount of edits the contributor has made to date.

TABLE 7:  
Regressions on the Relationship between Contributor Slant by Year and Prior Article Slant

Models	(1)	(2)	(3)	(4)	(5)	(6)
Dependent Variable	Contributor Slant by Year	Contributor Slant by Year	Contributor Slant by Year	Contributor Category by Year	Contributor Category by Year	Contributor Category by Year
Prior Article Slant	-0.0086*** [0.0001]	-0.0085*** [0.0001]	-0.0188*** [0.0004]			
Prior Article Category				-0.0147*** [0.0002]	-0.0147*** [0.0002]	-0.0228*** [0.0008]
Log(Prior Article Length)		0.0006*** [0.0000]	0.0010*** [0.0001]		0.0015*** [0.0000]	0.0021*** [0.0004]
Log(Prior Refs)		-0.0004*** [0.0001]	-0.0009*** [0.0001]		-0.0009*** [0.0000]	-0.0025*** [0.0004]
Observations	10,878,391	10,878,391	10,878,391	10,878,391	10,878,391	10,878,391
R-squared	0.006	0.007	0.007	0.001	0.001	0.001
Year FE	No	No	Yes	No	No	Yes
Article FE	No	No	Yes	No	No	Yes
Number of Articles	64,622	64,622	64,622	64,622	64,622	64,622

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

TABLE 8:  
Moderating Effect on How Contributor Slant Changes over the Years

Model	(1)	(2)	(3)	(4)
Dependent Variable	Abs(Contributor slant by year)	Abs(Contributor slant by year)	Abs(Contributor slant by year)	Abs(Contributor slant by year)
Average Bias of Articles Edited x Contributor Years	-0.0042*** [0.0001]	-0.0022*** [0.0004]		
Average Bias of Articles Edited	0.0174*** [0.0002]	0.0059*** [0.0008]		
Fraction of Extreme Articles Edited x Contributor Years			-0.0020*** [0.0001]	-0.0014*** [0.0004]
Fraction of Extreme Articles Edited			0.0088*** [0.0001]	0.0037*** [0.0006]
Contributor Years	0.0004*** [0.0000]	0.0001*** [0.0000]	-0.0001*** [0.0000]	-0.0001*** [0.0000]
Log(Number of Edits)	-0.0005*** [0.0000]	-0.0001*** [0.0000]	-0.0005*** [0.0140]	-0.0001*** [0.0000]
Observations	10,878,391	10,878,391	10,878,391	10,878,391
R-squared	0.011	0.011	0.009	0.008
Contributor FE	No	Yes	No	Yes

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

TABLE 9:  
Time Needed for a Contributor to Have >50% Probability of Moving to Neutral Slant

Starting Contributor Slant	Number of Years
Extremely Democratic	10
Democratic	6
Slightly Democratic	3
Neutral	0
Slightly Republican	4
Republican	7
Extremely Republican	11

Notes: Number of years calculated based on the Markov Chain Process. *Neutral* state includes contributor slant 0.5 standard deviation away from 0. *Slightly Democratic (Republican)* state includes contributor slant between 0.5 and 1.5 standard deviations below (above) 0. *Democratic (Republican)* state includes contributor slant between 1.5 and 2.5 standard deviations below (above) 0. *Extremely Democratic (Republican)* state includes contributor slant more than 2.5 standard deviations below (above) 0. On average, after about 30 years, the probabilities in all articles' end state reach stationary distribution, with the probability of contributor slant moving to *Neutral* being 87.4%.

TABLE 10:  
Heterogeneity of OA and BOF across Different Article Topics

Article Topics	No. of Edits	All sample		Republican contributors		Democratic contributors	
		Estimate	Pattern	Estimate	Pattern	Estimate	Pattern
Abortion	30,400	-0.0039*** [0.0012]	OA	-0.0161*** [0.0044]	OA	0.0003 [0.0012]	<i>n.s.</i>
Budget & Economy	765,729	-0.0019*** [0.0003]	OA	-0.0125*** [0.0011]	OA	0.0036*** [0.0003]	BOF
Civil Rights	902,531	-0.0038*** [0.0002]	OA	-0.0183*** [0.0008]	OA	0.0009*** [0.0002]	BOF
Corporations	54,709	-0.0009 [0.0008]	<i>n.s.</i>	0.0035 [0.0031]	<i>n.s.</i>	-0.0046*** [0.0007]	OA
Crime	957,613	-0.0016*** [0.0002]	OA	-0.0089*** [0.0009]	OA	0.0015*** [0.0003]	BOF
Drugs	164,330	-0.0029*** [0.0007]	OA	-0.0163*** [0.0025]	OA	0.0001 [0.0012]	<i>n.s.</i>
Education	864,373	-0.0064*** [0.0003]	OA	-0.0270*** [0.0011]	OA	-0.0028*** [0.0003]	OA
Energy	183,598	0.0021*** [0.0004]	BOF	0.0103*** [0.0015]	BOF	0.0012* [0.0007]	BOF
Family	434,980	-0.0013*** [0.0003]	OA	-0.0112*** [0.0014]	OA	0.0020*** [0.0004]	BOF
Foreign Policy	1,883,375	-0.0038*** [0.0002]	OA	-0.0079*** [0.0007]	OA	-0.0048*** [0.0004]	OA
Trade	442,561	-0.0038*** [0.0004]	OA	-0.0028*** [0.0010]	OA	-0.0125*** [0.0009]	OA
Government	3,376,993	-0.0039*** [0.0000]	OA	-0.0174*** [0.0004]	OA	-0.0026*** [0.0001]	OA
Gun	62,668	-0.0037*** [0.0009]	OA	-0.0207*** [0.0033]	OA	-0.0003 [0.0012]	<i>n.s.</i>
Healthcare	385,659	-0.0004 [0.0004]	<i>n.s.</i>	0.0027** [0.0014]	BOF	-0.0028*** [0.0006]	OA
Homeland Security	478,796	0.0021*** [0.0004]	BOF	0.0045*** [0.0014]	BOF	0.0025*** [0.0004]	BOF
Immigration	255,461	-0.0035*** [0.0005]	OA	-0.0031* [0.0019]	OA	-0.0047*** [0.0007]	OA
Infrastructure & Tech	920,016	-0.0017*** [0.0003]	OA	-0.0009 [0.0009]	<i>n.s.</i>	-0.0034*** [0.0004]	OA

Jobs	693,295	-0.0023*** [0.0003]	OA	-0.0074*** [0.0011]	OA	-0.0031*** [0.0004]	OA
Principles & Values	562,908	-0.0027*** [0.0003]	OA	-0.0017 [0.0012]	<i>n.s.</i>	-0.0071*** [0.0004]	OA
Social Security	2,501	-0.0111** [0.0048]	OA	-0.0365* [0.0190]	OA	-0.0138*** [0.0029]	OA
Tax	46,048	0.0058*** [0.0007]	BOF	0.0177*** [0.0033]	BOF	0.0039*** [0.0007]	BOF
War & Peace	1,837,644	-0.0018*** [0.0002]	OA	-0.0030*** [0.0007]	OA	-0.0022*** [0.0003]	OA
Welfare & Poverty	439,851	-0.0031*** [0.0004]	OA	-0.0109*** [0.0014]	OA	-0.0010** [0.0004]	OA
Biographies	1,311,337	-0.0024*** [0.0002]	OA	-0.0014* [0.0008]	OA	-0.0027*** [0.0003]	OA

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%; n.s.: not significant.

TABLE 11:  
Regression between Contributor Slant and Percentage of Republican in the Area

Model	(1)
Dependent Variable	Contributor Slant
RepPerc	0.0036*** [0.0013]
Observations	53,922
Adjusted R-squared	0.0001

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

TABLE 12:  
Relationship between Contributor Slant and Prior Article Slant, First Edits Only

Models	(1)	(2)
Dependent Variables	Contributor Slant	Contributor Slant
Prior Article Slant	-0.0092*** [0.0001]	-0.0218*** [0.0004]
Log(Prior Article Length)	0.0007*** [0.0000]	0.0011*** [0.0001]
Log(Prior Refs)	-0.0004*** [0.0000]	-0.0011*** [0.0001]
Observations	7,113,130	7,113,130
R-squared	0.007	0.007
Year FE	No	Yes
Article FE	No	Yes
Number of Articles	66,389	66,389

Notes: Robust standard errors in brackets. \*significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. Observations in this panel only include every contributor's first edit of an article.

FIG. 1. – Distribution of Contributors' Total Number of Edits

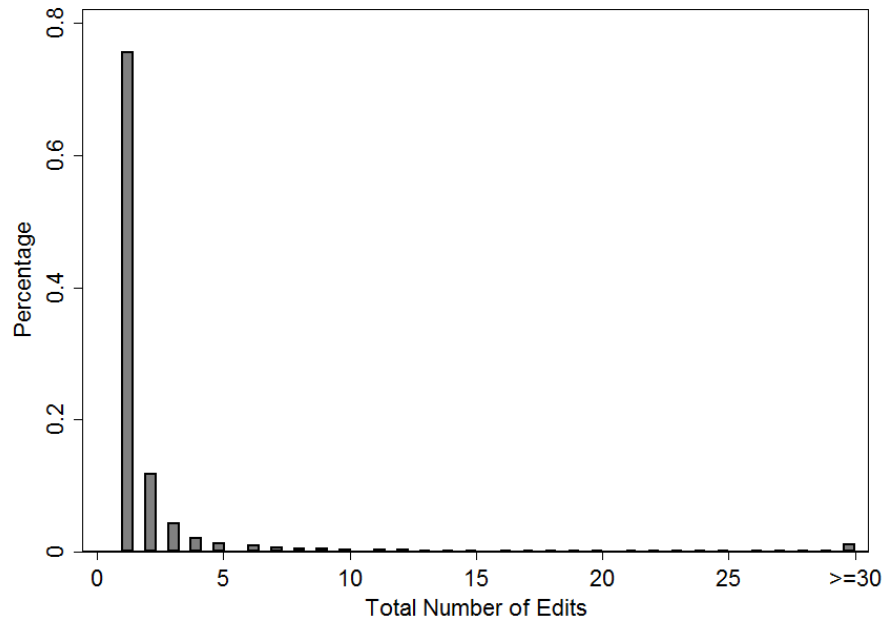


FIG. 2. – Transition Matrix of Contributor Slant Change in Wikipedia

		First half of activity		
		Democratic Type	Neutral	Republican Type
Second half of activity	Democratic Type	0.1407	0.0328	0.1145
	Neutral	0.7451	0.9333	0.7416
	Republican Type	0.1142	0.0339	0.1439

Notes: The sample is constructed by dividing every contributor's time in half. Then divide the direction of his or her edits, i.e. attach values (-1, 0, 1) to negative, 0, positive slant edits. Sum up the edits' values for the first half and the second half of his or her activity. If the sum of all edits in this half is negative, the contributor is a Democrat Type in this half. If the sum of all edits in this half is zero, the contributor is Neutral in this half. If the sum of all edits in this half is positive, the contributor is Republican Type in this half.



FIG. 3. – Transition Matrix of Contributor Slant Change over time

		Start							
		bin1	bin2	bin3	bin4	bin5	bin6	bin7	
Slant Range		[-1.229, -0.059)	[-0.059, -0.035)	[-0.035, -0.012)	[-0.012, 0.012)	[0.012, 0.035)	[0.035, 0.059)	[0.059, 1.000)	
End	bin1	[-1.229, -0.059)	0.8298	0.0139	0.0024	0.0011	0.0013	0.0008	0.0015
	bin2	[-0.059, -0.035)	0.0717	0.7242	0.0044	0.0020	0.0103	0.0019	0.0007
	bin3	[-0.035, -0.012)	0.0591	0.1745	0.7438	0.0055	0.0040	0.0149	0.0029
	bin4	[-0.012, 0.012)	0.0323	0.0713	0.2286	0.9795	0.2089	0.0531	0.0277
	bin5	[ 0.012, 0.035)	0.0036	0.0128	0.0177	0.0060	0.7545	0.1867	0.0624
	bin6	[ 0.035, 0.059)	0.0008	0.0014	0.0015	0.0033	0.0052	0.7222	0.0757
	bin7	[ 0.059, 1.000)	0.0028	0.0019	0.0018	0.0025	0.0158	0.0203	0.8291

Note: *Contributor Slant by Year* is split by the  $\pm 0.5$ ,  $\pm 1.5$ , and  $\pm 2.5$  standard deviations intervals. The middle bin represents neutral slant; the first/last bin represents extreme slant.

FIG. 4. – Vintage Analysis for Contributors Entering in Different Years

