# Battleground State:
# Characterizing and Understanding Political Bias in News Discourse on Reddit

**Mansi Ganatra**
University of Southern California
Computer Science Department
Los Angeles, CA 90089
mganatra@usc.edu

**Samar Haider**
University of Southern California
Computer Science Department
Los Angeles, CA 90089
samarhai@usc.edu

## Abstract

Increasing volumes of online content have made social media websites and news aggregators the primary medium of news readership for a majority of Internet users. While this offers a positive change in terms of the amount and variety of information available to people, there are certain downsides to such trends. Social media has been recently seen to amplify the biases that already exist in the real world. By allowing people from across the world to connect and partake in discussions with each other on a common platform, it creates an opportunity for communities to influence public forums and gain control over the conversation. This can lead to an over-representation of certain views, the creation of echo chambers, and more polarization in an already fractured society. The butterfly effects of such biases have resulted in significant sociopolitical changes in the world in recent years. In this paper we analyze the news discourse on Reddit, often called 'the front page of the internet', during the 2016 U.S. presidential election to uncover these bias and inform platform moderation policies to minimize their negative effects. Exploring activity and trends in neutral and partisan subreddits, we find different types of bias at the media, community, and linguistic levels.

## 1 Introduction

With more people making their way online every day, the Internet has arisen as the primary medium of information dissemination and large-scale human interaction. This has led to the creation of multiple platforms and social media websites that serve news and announcements from across the world. The increased spread of news, in terms of volume, velocity, and variety, is clearly beneficial to the world at large and makes for a more informed populace. But the functioning of online platforms is not like traditional news sources where the same version of a newspaper, curated by journalists, is distributed across the country. Online platforms, on the other hand, are primarily user-run, wherein participants themselves share news stories on a common discussion board for other users to engage with using their upvotes and comments.

This dynamic creates the opportunity for a number of different biases to creep into the online news readership model. Members of a community, and even those outside of it, can be biased in the topics of news that they choose to share as well as the sources they cite. Media outlets themselves have been shown to exhibit partisan bias, wherein some are more likely to cover left-leaning news in comparison to others, and vice versa. And finally, users themselves have the power to engage more strongly with certain posts over others, thereby boosting their ranking in the news feed and influencing the opinions and perceptions of other users on the website.

These biases do not fall under the umbrella of fairness in its traditional sense, which is driven primarily by acts of law that disallow unequal treatment of people along certain attributes. Yet these biases still capture an inherent inequality in the online world, one that may not be protected by law, but is still an important representation of society through the lens of social media. Biases that occur in the real world are amplified in the online space, and the widespread use of digital media for news dissemination has seen an increase in the formation of echo chambers. An imbalanced news diet and polarized political discourse stemming from it is detrimental to the public opinion and awareness. This in turn leads to the political landscape growing increasingly polarized over time, something the United States has witnessed happening in the Congress, resulting in a trend of hyper-partisan politics that has led to less cross-aisle cooperation than ever before. This is driven by the pressure on the Congress members by their constituents, who are driven by the news they consume and political discourse they engage in.

It is thus important to understand the biases that exist in political discourse and news readership on social media and aggregation platforms, especially on highly popular websites like Reddit, which have the power to influence millions of people. Most prior work in this direction has largely focused on Twitter as a source of data for political discourse centered around election campaigns. However, Twitter does not lend itself to the long-form discussions that take place on Reddit, which can often take the form of argumentation and play a role in swaying the opinions of people who might be on the fence on certain topics. On the other hand, work

that does look at activity on Reddit focuses more on the differences between certain communities and does not investigate the impact and engagement of partisan communities and users on neutral ground, which on its surface appear to be nonpartisan and unbiased.

In this project we investigate the nature of news readership and political discourse on the Internet in the age of social media and community-driven news aggregators. We use Reddit, the most popular social aggregation platform in the world, as a lens through which to understand the conversation surrounding the 2016 U.S. Presidential Election. We analyze the conversation on both neutral and partisan subreddits during the last three months of 2016 to uncover any latent biases in them. We find that there exist a number of such biases that offer insights into the thoughts of the online community during the election.

In the following sections, we focus on the following five research questions:

- **RQ 1:** Is there a media bias in content shared on news and political discourse subreddits? **Yes.**

- **RQ 2:** Does the identity of a media outlet reporting a news story impact a post's popularity and engagement levels? **Somewhat.**

- **RQ 3:** Is there a bias in the political discourse in news subreddits in terms of community engagement of users from partisan communities? **Yes.**

- **RQ 4:** Is there a difference in the sentiment of comments mentioning entities from both sides? **Yes.**

- **RQ 5:** Is there a geographical bias in the semantic representation of political entities in comments from local and global subreddits? **Yes.**

## 2    Related Work

A number of studies have focused on social media activity surrounding the 2016 United States presidential elections. (Bessi and Ferrara 2016) found that a significant fraction of the online conversation centered around the elections was driven by automated accounts. This is further expanded upon in (Badawy, Ferrara, and Lerman 2018), which uncovers the interference of the Russian IRA in the election, finding that conservatives were far more likely to be influenced by Russian bots than liberals. Similar analysis has been done on the 2017 French elections (Ferrara 2017) and the 2017 German elections (Morstatter et al. 2018). Prior work on Reddit has focused mostly on its effectiveness and fairness as an aggregator as in (Gilbert 2013), and the impact of anonymity on the nature of online conversations (De Choudhury and De 2014). (Singer et al. 2014) find that, while topics are becoming more diverse on Reddit, there has emerged a concentration towards certain domains as represented by popularity and engagement. (Kumar et al. 2018) show the prevalence of echo chambers in inter-community conflicts on Reddit, and the issue of bias in popularity of content is analyzed in (Lakkaraju, McAuley, and Leskovec 2013), which shows that a number of factors predict how communities engage with new submissions.

## 3    Data

We use PushShift[1] as our source of data. We download the dumps of submissions and comments for October, November, and December 2016, which totals over 160GB. We filter them to only retain data from our 9 target subreddits: 3 neutral, 3 left-leaning, and 3 right-leaning. This gives us more than 1 million submissions and 16 million comments. Table 1 shows the number of submissions, comments, and file size of the data for each of the subreddits.

| Subreddit | # of submissions | # of comments |
|---|---|---|
| r/news | 199135 | 1798331 |
| r/worldnews | 101965 | 1989895 |
| r/politics | 126173 | 6700996 |
| r/hillaryclinton | 22608 | 441156 |
| r/The_Donald | 674119 | 5120822 |
| r/democrats | 5918 | 17231 |
| r/Republican | 2286 | 19388 |
| r/Liberal | 2125 | 7297 |
| r/Conservative | 11660 | 131381 |
| Total | 1,145,989 | 16,226,497 |

Table 1: Number of submissions and comments in each of our target subreddits during October, November, and December 2016. We see that the r/politics subreddit has the highest number of comments and comments-to-submissions, which shows that it is primarily used as a discussion and debate forum.

## 4    Media Bias

In this section, we look at the difference in the number of submissions as well as their popularity and engagement for news media with different political leanings. We classify submissions as quoting either left-, center-, or right-leaning news sources as defined by AllSides[2], an independent news media rating organization. This gives us a list of 47 news websites in total. Table 2 shows the number of submissions quoting websites which lean toward different political ideologies.

| Subreddit | Left | Center | Right | Social media |
|---|---|---|---|---|
| r/news | 10604 | 7583 | 2786 | 29307 |
| r/worldnews | 8112 | 10386 | 1562 | 14397 |
| r/politics | 32582 | 10720 | 6947 | 15783 |
| Total | 51298 | 28689 | 11295 | 59487 |

Table 2: Number of matches found for news websites with different political leanings as defined by AllSides. Left-leaning websites are generally over-represented in comparison to right-leaning ones.

### 4.1    Comparison of Media Bias in Subreddits

Figure 1 shows the bias in the popularity of and engagement on submissions which quote news websites with different political stances. We see a significant left-leaning bias in

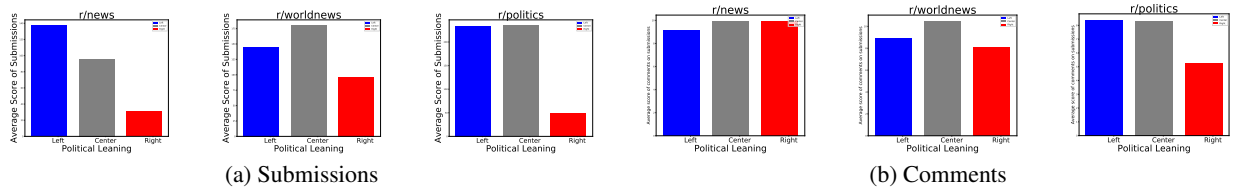(a) Submissions



(b) Comments

Figure 1: A comparison of the bias in engagement levels for links from left-, right-, and center-leaning news websites. Subfigure 1a shows the average score of submissions and Subfigure 1b shows the average score of comments on those submissions by political leaning in each of the three news subreddits.

submissions in the U.S.-centric r/news. The global r/world-news is more centrist in comparison, while r/politics shows a more centre-left bias. In terms of comments, there is less of a difference, which shows that the engagement is somewhat balanced.

## 4.2 Analysis of Submission Popularity

We investigate whether a submission's popularity is driven by the media bias of the news website mentioned in it. We create a dataset with features including semantic representation of the title, the political leaning of the news source, and the time of day that the submission is made. We measure the popularity of a submission in terms of the number of upvotes and comments it gains. Table 3 shows the correlation between the features and the number of comments and upvotes on the submission. We see that left-leaning media bias shows a higher correlation with popularity than other features. We use analysis of variance (ANOVA) to compute the significance of these results.

| Feature | Value | # of comments | # of upvotes |
|---|---|---|---|
| Media bias | Center | 0.0196 | 0.01528 |
| | Left | **0.04158** | **0.03076** |
| | Right | 0.00541 | -0.00072 |
| | Other | -0.0177 | -0.01026 |
| *p-value* | Overall | $1.20e^{-100}$ | $1.70e^{-56}$ |
| Time | Morning | 0.01254 | 0.01148 |
| | Noon | 0.01138 | 0.01086 |
| | Evening | -0.01392 | -0.01322 |
| | Night | -0.00727 | -0.00657 |
| *p-value* | Overall | $1.14e^{-17}$ | $1.54e^{-15}$ |

Table 3: The correlation between the features and the number of comments and upvotes on a submission on r/news. We see that the political leaning can help predict the eventual popularity of a submission.

## 4.3 Temporal Trends in Bias during Elections

Figure 2 shows the media bias trends over time in all three news subreddits during election season. We see left-leaning dominate right- and center-leaning ones, and a spike in submissions quoting social media sources in the day of the election, November 8, which can be attributed to the surge in tweets after Donald Trump won.

## 5 Bias in Partisan Community Engagement

One of the biggest issues in online political communities is the formation of echo chambers which leads to confirmation bias; showing people views and opinions that match their own leads them to think that their beliefs are more prevalent than they actually are. This is one of the reasons Donald Trump's victory came as a shock to many liberals. They were so entrenched in echo chambers where everyone held the same leftist views that they did not consider the vast number of people outside their community—who were themselves stuck in their own echo chambers. This phenomenon also caused increased polarization in political views; with no one to oppose or offer a different opinion, people tend to believe in their own perspective even more strongly. The widening chasm between Democrats and Republicans in recent years can be attributed to this phenomenon.

One way to break out of these echo chambers is to read and partake in discussions in neutral communities, such as the news and politics subreddits. However, if one political party is more engaged in the discussion than the other, it can lead to the creation echo chambers similar to those that already exist in partisan communities. Furthermore, it may give the impression to the casual reader that that a certain kind of political opinion is the norm on Reddit.

Here we compare the bias in engagement levels of left- and right-leaning accounts from the neutral community's perspective. We define cross-posters as people who have posted in both a partisan community as well as a neutral one during election. The definition of partisan community engagement in neutral subreddits is as follows:

$$E_{S \to T} = \frac{|\{u : u \in S \text{ and } u \in T\}|}{|T|} \quad (1)$$

where $E_{S \to T}$ is the engagement level of users from the source subreddit in the target subreddit, $u$ is a unique user account, $S$ is the set of authors from the source subreddit (in this case the partisan community), $T$ is the set of authors from the target subreddit (in this case the neutral community).

Figure 3 shows the relative difference in engagement from right- and left-leaning communities in the news subreddits. We see that users from right-leaning subreddits are over-represented in comparison to those from left-leaning subreddits as a ratio of the total users in news subreddits. This can influence neutral readers into thinking that right-wing views are the norm on Reddit.
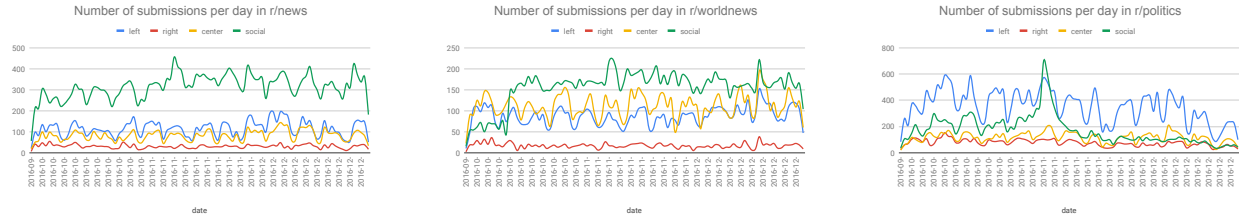
Figure 2: Media bias trends over time in all three news subreddits during election season. References to social media sees a spike on November 8, the day of the election.
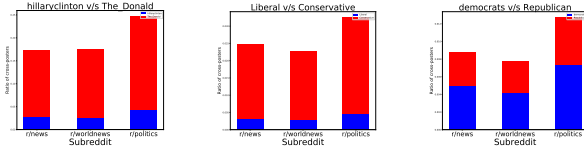


Figure 3: A comparison of the bias in levels of engagement by users from left- and right-leaning communities in the r/news subreddit. The chart depicts the ratio of users active in both subreddits and the total number of users who post in r/news. We see that users from r/The_Donald and r/Conservative far outnumber users from r/hillaryclinton and r/Liberal, respectively.

## 6 Linguistic Bias in Political Discourse

In this section we look at the difference in language used when referring to political actors and entities in local and global news subreddits, as well as the contrast in linguistic patterns between left- and right-leaning subreddits.

### 6.1 Sentiment Analysis of Comments

We extract all comments that contain mentions of the major political actors and entities involved in the election, namely Hillary Clinton, Donald Trump, the Republican and Democratic parties, and the terms liberal and conservative. We then evaluate the sentiment expressed in those comments using VADER (Hutto and Gilbert 2014), a tool for calculating the sentiment of texts on social media using a rule-based lexicon. We compare the average sentiment in r/news and r/worldnews to understand how the the local U.S. population and the rest of the world perceive the two candidates. Figure 4 shows the comparison between the sentiment surrounding the different keywords in local and global subreddits. We can see that overall, left-wing candidates terms have a more negative connotation attached to them than right-wing ones. It is surprising to see that comments in both domestic and international communities are more favorable towards Donald Trump and the Republican party than they are towards Hillary Clinton and the Democractic party. However, this result is corroborated by (Ferrara 2017), who show that a similar trend was found to exist on Twitter during the 2016 elections. While they found that this trend was driven by bots on Twitter, we see here that a similar exists on Reddit, which might prompt the question about the prevalence of similar bot accounts on Reddit.
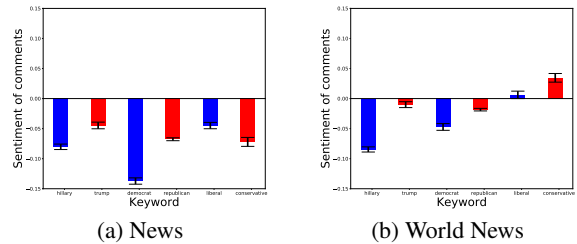


(a) News      (b) World News

Figure 4: A comparison between sentiment scores for comments containing left- and right-leaning keywords on r/news and r/worldnews. We can see that in both r/news and r/worldnews, a stronger negative sentiment is expressed for Hillary Clinton and the Democratic party as compared to Donald Trump and the Republican party, respectively. We include error bars to show that the results are significant.

### 6.2 Semantic Representations of Political Entities

Beyond sentiment, the semantic representation of political actors and entities can also provide clues about the biases that exist in the discourse in different communities. We first train word embeddings using Word2Vec (Mikolov et al. 2013) separately on comments from the r/news and r/worldnews subreddits. Table 4 shows the top 15 nearest neighbors in the vector space for the input words 'hillary' and 'trump' in both subreddits. We see that the local discussion is mostly centered around the relation between the two candidates and other local politicians and parties while the global discussion talks about their relation with other world leaders.

The Word Embedding Association Test (WEAT) (Caliskan, Bryson, and Narayanan 2017) is used to find biases in word associations between target entities and a series of attributes. Using a seed list of five attributes and their opposites (good/bad, smart/stupid, friendly/unfriendly, honest/dishonest), we use the English WordNet database (Miller 1995) to obtain synonyms for each of them, giving us an extensive list of 60 positive and negative adjectives. We train word embeddings separately on comments from communities with left (r/hillaryclinton, r/democrats, r/Liberal), right (r/The_Donald, r/republican, r/Conservative), and neutral (r/news, r/worldnews) political stance. We investigate the effect size between Hillary and Trump, Democrat and Republican, and Liberal and Conservative as defined in the original paper and compare the values obtained for all three models. The effect size compares the difference in attribute associations between different targets (e.g. Hillary vs. Trump) along the positive dimension, i.e. a

| r/news | | r/worldnews | |
|---|---|---|---|
| hillary | trump | hillary | trump |
| clinton | hillary | clinton | hillary |
| hilary | hrc | hilary | obama |
| hrc | hilary | hrc | putin |
| bernie | drumpf | bernie | hilary |
| trump | obama | sanders | clinton |
| sanders | trumpler | trump | drumpf |
| shillary | clinton | obama | duterte |
| stein | bernie | clintons | hrc |
| obama | djt | dnc | he |
| dnc | tump | gop | bernie |
| killary | sanders | killary | netanyahu |
| clintons | ukip | romney | sarkozy |
| hillarys | trumps | dws | bibi |
| romney | killary | shillary | du30 |

Table 4: The top 15 words most similar to 'hillary' and 'trump' in r/news and r/worldnews. We see that the international conversation is centered around the relation of the two candidates to other heads of state, while the domestic conversation primarily focuses on their local rivals.

higher score indicates a more positive association with the target, and is defined as follows:

$$Eff = \frac{mean_{x \in X} s(x, A, B) - mean_{y \in Y} s(y, A, B)}{std_{w \in X \cup Y} s(w, A, B)} \quad (2)$$

$$s(w, A, B) = mean_{a \in A} cos(\overrightarrow{w}, \overrightarrow{a}) - mean_{b \in B} cos(\overrightarrow{w}, \overrightarrow{b}) \quad (3)$$

Table 5 shows this comparison. We can see that Trump is more closely associated with positive attributes in right-leaning subreddits in comparison to Hillary, while the reverse is true for left-leaning ones. The same is true for other target words.

| | Community stance | | |
|---|---|---|---|
| Candidate | Left | Right | Neutral |
| Hillary vs. Trump | **0.1999** | -0.5512 | -0.0237 |
| Democrat vs. Republican | **0.1587** | -0.3939 | -0.0279 |
| Liberal vs. Conservative | -0.1128 | -0.6201 | **-0.1030** |

Table 5: Bias in WEAT effect sizes for positive attributes for Hillary vs Trump, Democrat vs Republican, and Liberal vs Conservative across different communities.

## 7 Conclusion and Future Work

Bias in social media and online communities is an important problem that can damage political discourse on the Internet and lead to real-world consequences. In this paper we analyzed the news discourse on Reddit during the 2016 U.S. presidential election and find a series of biases of different types that influenced the political conversation around it. In the future we plan to extend this work by conducting a deeper analysis of the trends we found and attempt to find reasons for their existence. We also plan to use model

the communities as networks to discover other websites with media bias similar to those listed by AllSides, which can be used as a more exhaustive list of sources of biased news reports on the Internet.

## References

Badawy, A.; Ferrara, E.; and Lerman, K. 2018. Analyzing the digital traces of political manipulation: the 2016 russian interference twitter campaign. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 258–265. IEEE.

Bessi, A., and Ferrara, E. 2016. Social bots distort the 2016 us presidential election online discussion. *First Monday* 21(11-7).

Caliskan, A.; Bryson, J. J.; and Narayanan, A. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356(6334):183–186.

De Choudhury, M., and De, S. 2014. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. In *Eighth International AAAI Conference on Weblogs and Social Media*.

Ferrara, E. 2017. Disinformation and social bot operations in the run up to the 2017 french presidential election. *First Monday* 22(8).

Gilbert, E. 2013. Widespread underprovision on reddit. In *Proceedings of the 2013 conference on Computer supported cooperative work*, 803–808. ACM.

Hutto, C. J., and Gilbert, E. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth international AAAI conference on weblogs and social media*.

Kumar, S.; Hamilton, W. L.; Leskovec, J.; and Jurafsky, D. 2018. Community interaction and conflict on the web. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*, 933–943. International World Wide Web Conferences Steering Committee.

Lakkaraju, H.; McAuley, J.; and Leskovec, J. 2013. What's in a name? understanding the interplay between titles, content, and communities in social media. In *Seventh International AAAI Conference on Weblogs and Social Media*.

Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 3111–3119.

Miller, G. A. 1995. Wordnet: a lexical database for english. *Communications of the ACM* 38(11):39–41.

Morstatter, F.; Shao, Y.; Galstyan, A.; and Karunasekera, S. 2018. From alt-right to alt-rechts: Twitter analysis of the 2017 german federal election. In *Companion Proceedings of the The Web Conference 2018*, 621–628. International World Wide Web Conferences Steering Committee.

Singer, P.; Flöck, F.; Meinhart, C.; Zeitfogel, E.; and Strohmaier, M. 2014. Evolution of reddit: from the front page of the internet to a self-referential community? In *Proceedings of the 23rd international conference on world wide web*, 517–522. ACM.