

Detecting Linguistic Change in Reddit's *The_Donald*

ADARSH MATHEW

Final Paper, Computational Content Analysis (Winter 2020)

adarshm@uchicago.edu

I. BACKGROUND

As the world increasingly comes online, one should expect a version of our political activities being performed on these platforms too. The affordances of the internet and networked platforms like Twitter, Reddit, and Facebook reduce the transaction costs of engagement at scale. By reducing the cost of content creation and downgrading the importance of traditional mediators of attention and political influence like the news media, platforms have enabled the amplification and shaping of erstwhile niche or marginalized narratives. These very affordances make political communication and activism a natural fit for actors on these platforms.

Political communication has always relied on tapping into pop culture to make its message more palatable to the public and likely voters. This has extended to digital politics too: political memes perform the role of satire, candidates maintain a presence across multiple platforms, which pollsters and campaigns rely on for fundraising and electoral analysis. But digital politics also means that a candidate or party's messaging on these platforms cannot be a controlled and mediated spectacle like it could have been in traditional media. Near-instant feedback and engagement, algorithmic amplification, and the sheer scale of platforms create an ecosystem where no campaign can reasonably hope to exert top-down control on narratives, messaging, and language.

If President Obama's re-election in 2012 was the first US presidential election that relied on social media for targeted messaging, Donald Trump's election win in 2016 marked the next stage of digital politics, driven by micro-targeted advertising, candidate tweets, data breaches, and disinformation campaigns. In addition to their messaging success, Trump's base proved to be extremely good at setting the agenda – re-framing the controversy around Hillary Clinton's e-mails as a referendum on her lack of honesty, the overtly racist framing on issues like 'Black Lives Matter', etc. Users of 4chan/8chan's r/pol/ and reddit's r/The_Donald were able to command attention online, driving their narrative through memes and highly specific language familiar to digital natives. In light of all this, how does one track the changing context around issues in these communities? How does a user's linguistic behaviour change with greater engagement? How does one measure the embedded context around political entities on these forums?

President Donald Trump's rise has been associated with the rise to prominence of far-right communities, branded as the 'alt-right'. The violent, ostensibly supremacist 'Unite the Right' rally in Charlottesville, Virginia was seen by many as a watershed moment that marked the rise of this subgroup (Atkinson, 2018). President Trump's election campaign itself was marked by several racist and misogynistic statements, which were cheered on wildly and amplified by this subgroup. Heikkilä argues that by actively antagonizing both Democrats and classic conservatives – through memes and online culture wars – the alt-right pushed the conversation farther to the right, normalized the vilification of 'political correctness', and created a feedback loop that rewarded the campaign for its extreme positions (Heikkilä, 2017; Nagle, 2017).

The tactics of the alt-right are heavily influenced by the 'ironic' meme culture of 4chan/8chan. The origins and inspirations for its more virulent side, which includes shitposting, bullying, targeted harassment, doxxing (release of private information), can be traced to 'GamerGate'. 'GamerGate' refers to the targeted harassment of video game journalists orchestrated through the hashtag '#GamerGate'. Beginning in August 2014, it involved the large-scale coordinated harassment of specific women video game journalists who were critical of the portrayal of women in video games. Characterized as an effort to resist the imposition of politically correct norms on the art they love, the community lashed out at specific targets. Victims received death threats, rape threats, their addresses were posted online (doxxing). GamerGate was a vital moment in the crystallization of the alt-right as we know it (Blodgett and Salter, 2018). The alt-right embraced the tactics of GamerGate, employing and co-opting the meme-first approach to discourse, along with a full-throated disavowal of feminist, politically-correct discourse. All of this was interspersed with 'ironic' jokes replete with racist, anti-Semitic, and misogynistic themes.

The importance of GamerGate in understanding the alt-right movement cannot be understated – the incident tapped into a well of discontent in a population that wielded a lot of power in the online world, but felt marginalized. The misogynistic nature of the alt-right's fundamental beliefs have been touched upon by academics and journalists. The alt-right is defined by its white supremacist and anti-Semitic ideology, but when alt-right members talk about the women who are standing between them and their "rightful" position, their language is virtually indistinguishable from what you can find on misogynistic message boards (League, 2018). The subculture is characterized by its extreme disavowal of 'political correctness' and its norms, along with the 'Social Justice Warriors' who enforce it online. They're strongly opposed to Feminist values; 'Feminism is Cancer' continues to be one of the most popular memes of the alt-right, and still serves as a rallying cry. Casting themselves as noble warriors defending the sanctity of what they love, young, white men – railing against the enforced 'politically correct' critiques of the video-game industry – engaged in coordinated, large-scale, targeted harassment while being cheered on gleefully (Nieborg and Foxman, 2018). The actions of key players set the norms for the community, and normalized a vicious version of internet 'trolling', giving birth to a plethora of celebrities like Milo Yiannopoulos, who modelled himself as an 'anti-feminist, ultra-conservative bad boy' (Koulouris, 2018). It provided a blueprint for reactionary mobilization and engagement on the internet (Daniels, 2018).

In this paper, I attempt to understand how linguistic behaviour of users in online communities have changed over time, and what the relationship between the language of prominent users and the community at large looks like. I also seek to quantify how the context around key phrases has changed over time using word embedding methods. Additionally, using word embeddings, I construct dimensions of key electoral issues and track how the perception of political figures and institutions shift on these dimensions.

II. METHODS

Data Sources

I use data from the *r/The_Donald* subreddit to track and understand changes in user linguistic and rhetoric behaviour over time. Using a combination of the Reddit API¹ and the pushshift.io API² (Baumgartner et al., 2020), I use all submissions and comments made from 16 June 2015 to 8 November 2018 as my corpus of text. This totals almost 2.3 million unique entries for a period of 40 months. The data collection tracks the subreddit from when President Trump announced his candidacy to the two-year anniversary of his election win.

Snapshot Language Models & User Change

To quantify changes in user linguistic behaviour relative to the community, I adapt the methodology used by Danescu-Niculescu-Mizil et al.’s 2013 paper *No Country for Old Members* (Danescu-Niculescu-Mizil et al., 2013). I use a bi-gram Snapshot Language Model (*SLM*) with normalized tokens. I use the Interpolated Kneser-Ney smoothing method, which mixes a discounted probability with a lower-order continuation probability, with a discount factor of $d = 0.1$ ³. Kneser-Ney smoothing makes use of the probability of a word being a novel continuation. While Danescu-Niculescu-Mizil et al use a Katz back-off smoothing approach tuned on a held-out set, I choose Kneser-Ney Interpolation for its superior performance (Jurafsky and Martin, 2008). The duration of each snapshot is 6 weeks, thus generating 30 such windows to examine user and sub-group behaviour.

For each post in a six-week window, I compute its cross-entropy according to the $SLM_{w(p)}$ of the window. This cross-entropy is a measure of how surprising the language of the post is relative to the language of the subreddit in the window. Higher cross-entropy values indicate greater deviation from the language of the subreddit. Like Danescu-Niculescu-Mizil et al., I only consider the first 5 sentences of each post to calculate the cross-entropy metric; longer posts would have higher cross-entropy, so a hard limit of $k = 5$ makes it possible to compare cross-entropy scores across posts. Due to considerations of computation time, I only examine submissions and not comments for this method. Additionally, I calculate cross-entropy for posts which have 10

¹Reddit API: <https://www.reddit.com/dev/api/>

²Pushshift API: <https://pushshift.io/api-parameters/>

³The default discount factor in Python’s `nltk` library is set to 0.1.

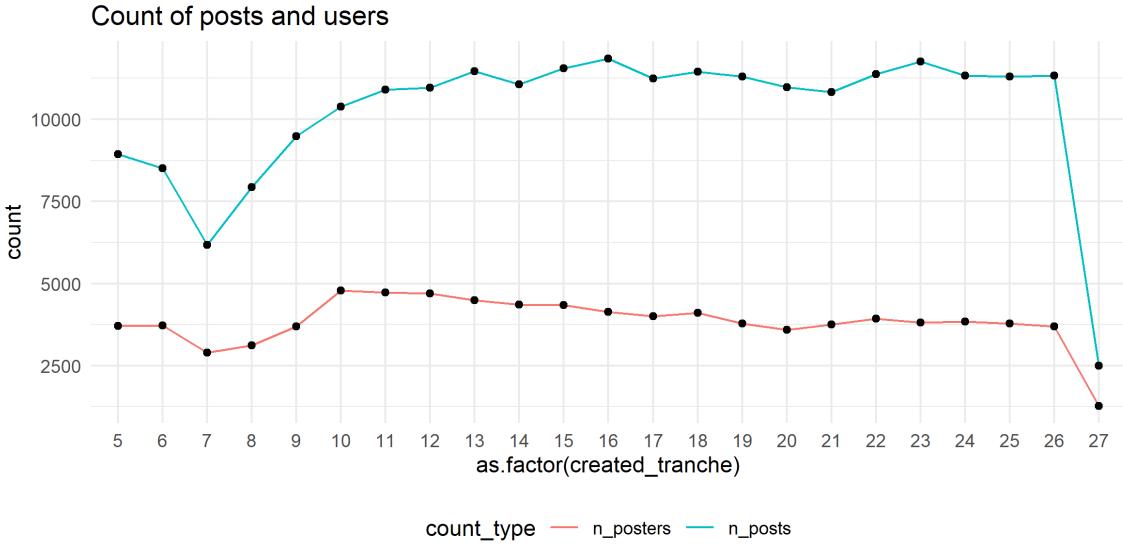


Figure 1: Count of posts and posters for Snapshot Language Models.

comments or more, have a score of 100 or more. I drop the first 5 and last period, giving us almost 230k posts over 23 periods to examine.

Word Embeddings to describe Context & Change in Community Discussions

To describe the context surrounding topics of interest in the community, I construct word embedding models for each period using skip-grams with negative sampling using Python’s `gensim`(Řehůřek and Sojka, 2010) library. The high-dimensional corpus of text is reduced to lower dimension where each word w_i is represented by a vector \mathbf{w}_i of size k (Kozłowski et al., 2019). The skip-gram approach captures the nature of co-occurrence between words and phrases. Semantic similarity – and divergence, by extension – between two words is measured using the cosine similarity between the two word vectors representing them (Hamilton et al., 2018). I retain all posts and comments for this method.

To compare the divergence of words from different time-periods, we need to ensure that the word vectors are aligned to the same coordinate axes. I use the Procrustes alignment method developed by Hamilton et al. (Hamilton et al., 2018). For p time-periods, we take the common vocabulary across these periods and obtain the optimal rotational alignment that preserves cosine similarities. After alignment, we can use the adjusted word vectors to compute semantic divergence, i.e. how the context of a word has shifted over time. This is measured using the cosine distance between the word vectors in time $t, t + \Delta$: $\text{cosineDist}(\mathbf{w}_t, \mathbf{w}_{t+\Delta})$. (Hamilton et al., 2018)

Additionally, I construct semantic dimensions for issues to track how certain the context surrounding specific issues shift over time. For each dimension, I take the normalized sum of antonym pairs to create a bi-directional vector to project word vectors onto. (Kozłowski et al., 2019) The cosine-similarity of word vectors with these dimensions tells us where the word lies on

the spectrum created by that dimension.

I create six dimensions of interest, which are enumerated below. The complete list of antonym pairs used to construct them is provided in Table 1.

1. Security: How closely do entities relate to issues of safety and security. Eg: Are candidates *strong* or *weak* on crime?
2. Trade: To capture the context around trade issues – *local* protectionism vs. *global* free trade.
3. Ideology: Create a dimension to quantify the ideological spectrum between *left* and *right*, *conservative* and *liberal*.
4. Region: Are the issues of the party focused on cities and the *coast*, or are they referring to the *heartland*?
5. Economy: Is the party associated with the *rich* and the elite, or do they speak to the *poor* in America.
6. Race: Where are candidates positioned on questions of race and immigration.

Additionally, I restrict my analysis of these semantic dimensions to four distinct windows. This is to simplify the attempt to explain sifts within the same dimension driven by events. These windows are selected to indicate vital time periods in the Trump presidency. These periods include:

1. Super Tuesday, March 2016: A major victory for Donald Trump in the Republican Primary in key states on Super Tuesday, which made him the front-runner for the nomination. In the run-up to Super Tuesday, he moved from being an outside candidate to a viable candidate, fleshing out his policies and campaign rhetoric.
2. General Elections, November 2016: Donald Trump beats Hillary Clinton in an upset victory by the slimmest of margins. Coverage in the final days was heavily influenced by investigations into Clinton's e-mail server, precipitated by a letter from FBI director James Comey to Congress about newly discovered material days before the election.
3. The first quarter of the presidency, February - March 2017.
4. 'Unite the Right' rally in Charlottesville, Virginia, August 2017: This was a white-supremacist rally conducted in response to the removal of Confederate monuments by local governments. The rally served as a unifying event for the American white-nationalist movement, and attracted counter-protestors. It turned violent after clashes between the two groups, and on Day 2, a white supremacist deliberately drove his car into the counter-protestors, killing one of them.

Table 1: Antonym pairs used to construct Semantic Dimensions

Security	Trade	Ideology	Region	Economy	Race
safe - vulnerable	local - global	right - left	rural - urban	rich - poor	white - black
strong - weak	tariff - free	conservative - liberal	town - city	advantage - need	american - immigrant
			interior - coast	equity - loan	close - open

III. RESULTS

Surprising Language and User Behavior

The distribution of cross-entropy for each period is fairly similar – we notice each distribution is fat-tailed and the mean cross-entropy is stable. Period 10, the month of the election, has the highest average cross-entropy. There does not seem to be any variation with time. From Figure 3, we also notice that there is no relationship between the cross-entropy and score of a post, and that the relationship between cross-entropy and the number of comments is weak. Thus, filtering on these two parameters would not be equivalent to selecting on the outcome.

To understand the influence of top users, I plot the average cross-entropy of their posts against that of the community in Figure 4. Top users are defined as those who have made at least 5 posts per period, and whose posts have a mean score greater than 2500 – regular users whose posts get consistently high scores, about the top 3% of posters and 10% of posts. We notice that the cross-entropy of these top users is highly correlated to that of the community at every period, indicating they do not have wildly different language usage from the community. Additionally, we observe that the average cross-entropy of these users is less than the average of the community, i.e. their language is less surprising. This could be an indicator that these users embody the norms of the community, and are thus able to command attention.

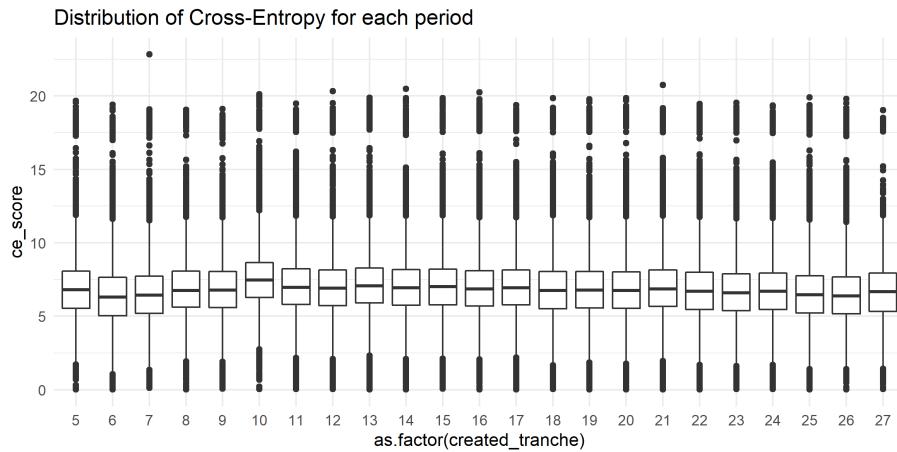


Figure 2: Distribution of Cross-Entropy for each 6-week period.

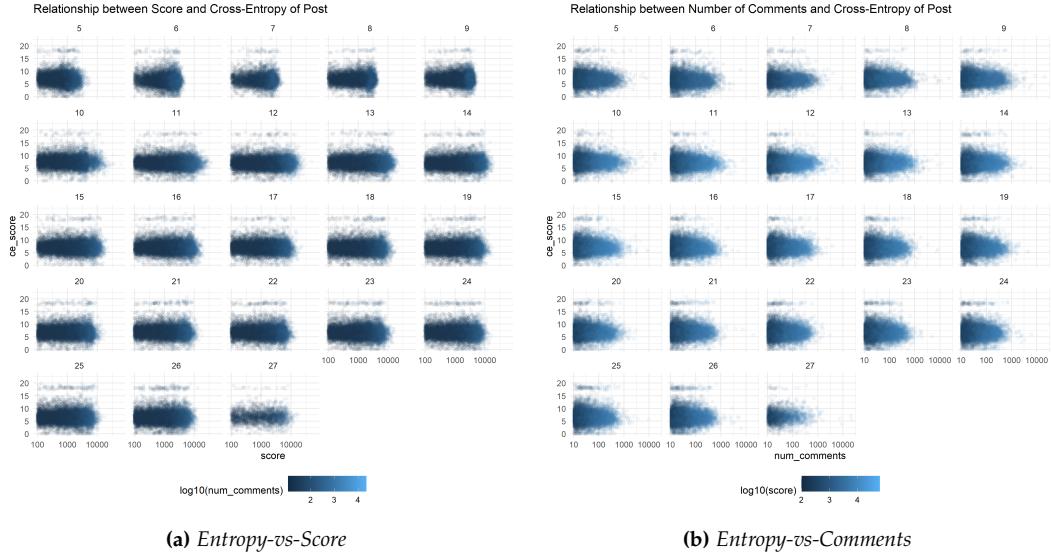


Figure 3: Relationship between Cross-Entropy and Score, Number of Comments.

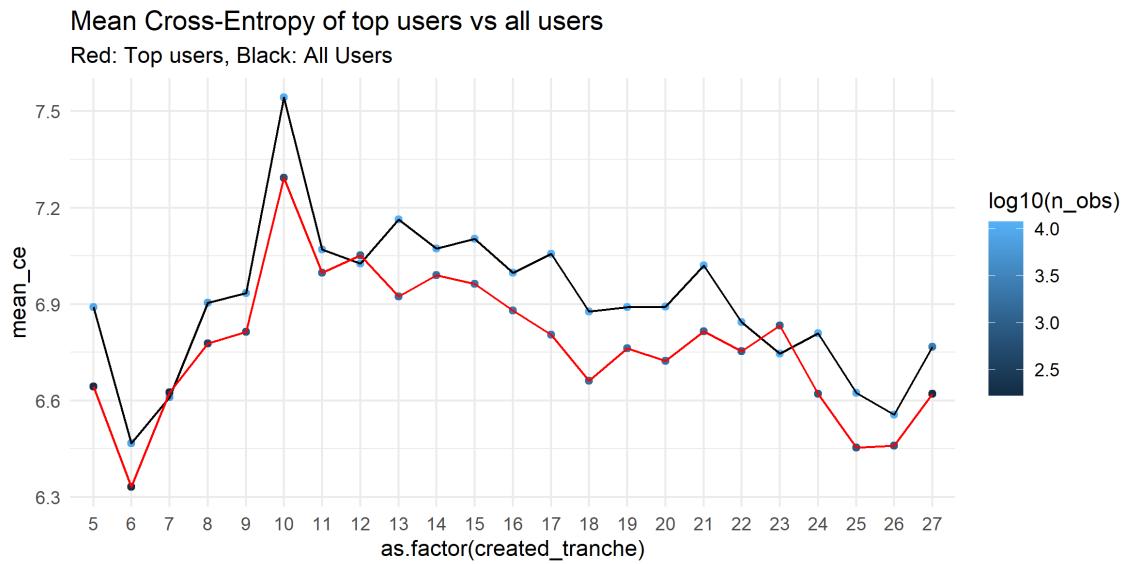


Figure 4: Mean Cross-Entropy of all posts vs top users for each 6-week period.

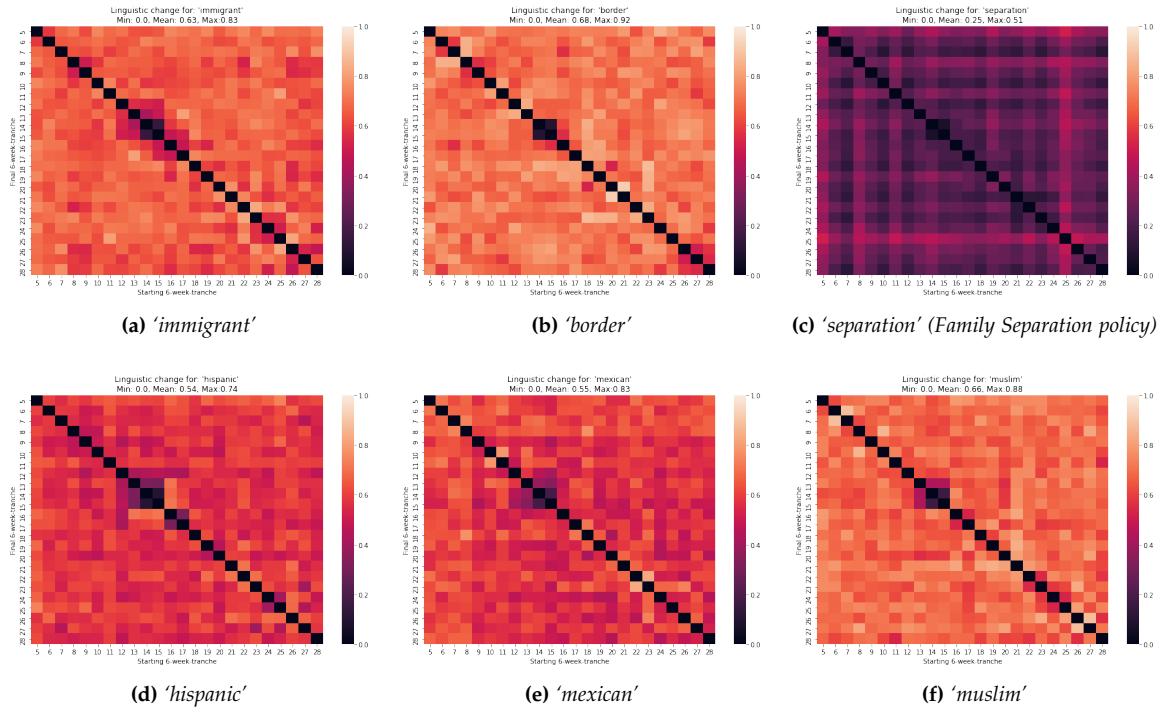


Figure 5: Relative Divergence of words related to Trump’s Immigration policies using aligned word embeddings.

Using Divergence measures to identify a change in context

Immigration policies

Immigration was and continues to remain a cornerstone of President Trump’s legislative agenda. The first notes of stabilization for `immigrant`, `border` begin in period 12 in the early months of the presidency, as we see in Figure 5. The thrust on immigration continues until period 17, before it picks up again during periods 22. This later spike is associated with the legislative focus on the border family separation issue. The community’s context around `separation` changes in period 25 – when President Trump signed the executive order rescinding the Family Separation policy (`noa`).

The focus on immigration translates into increased attention on ethnic and religious minorities. There is a notable change in the use of `mexican` right after the election, with an associated and muted change in `hispanic` a few periods later, with both stabilizing in the same periods as above. I expected to see a related uptick in `muslim` in association to President Trump’s Muslim ban (`noa`, 2017). The ethno-nationalist focus of the subreddit’s ideology creates confounders in contextualizing the word.

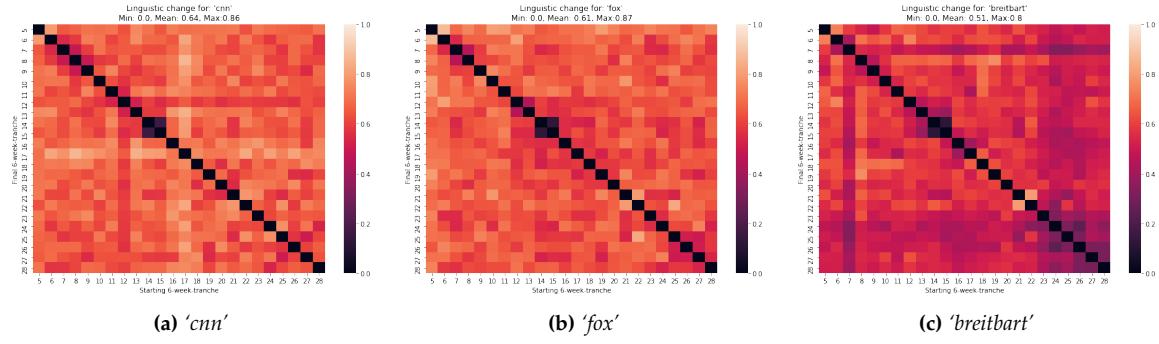


Figure 6: Relative Divergence of names of prominent news outlets using aligned word embeddings.

References to news media

A salient highlight of President Trump’s term has been his frequent attacks on news media outlets, charging them as the ‘enemy of the people’, of having a vendetta against him, going so far as to institute a list of ‘Fake News’ award winners. This messaging has had significant success within his base, furthering a partisan re-framing of institutions and working to de-legitimize them. The multiple and varying references to *cnn* – albeit in different contexts, since they report on a multitude of issues – seem like a strong indicator of their centrality to the conversations in the subreddit. The spike in divergence values in period 17 refer to the incident where President Trump re-tweeted a video of attacking CNN.

On the other hand, President Trump has also expertly tapped into the right-wing media ecosystem to amplify his message. He has made several announced appearances on talk shows hosted by Fox News to chime in on policy. Steve Bannon, chief editor of Breitbart News, was a key figure in the election campaign and was appointed Chief Strategist after the election win, at par with the Chief of Staff (noa, 2016). Breitbart has also been a cornerstone of the alt-right, with its nativist and highly partisan coverage crafted for online amplification, replete with fake news and conspiracy theories. For two important news outlets occupying a similar ideological space, their impressions in the subreddit vary drastically.

fox’s relative divergence map varies in a form similar to that of *cnn*, while *breitbart* is likely used in fewer contexts. For the latter, we notice a high change in context at period 22, followed by a period of lower divergence, indicating a stabilization and narrowing of context. The relative fortunes of Fox and Breitbart indicate a shift in attention of the subreddit’s sources, with Fox taking on a markedly different role from the first month of the presidency. This shift away from Breitbart is captured well in the subreddit’s references to *milo*, i.e. Milo Yiannopoulos, alt-right provocateur and former editor of Breitbart News. Yiannopoulos’s decline as the darling of the alt-right – following his remarks on pedophilia and child sexual abuse – tracks with the muted attention of Breitbart in the subreddit, an indicator of weakening influence.

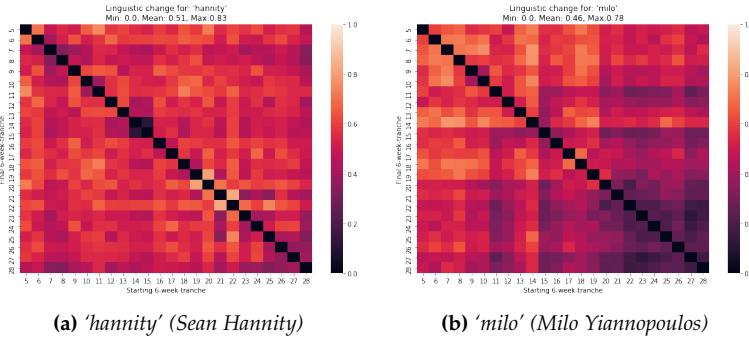


Figure 7: Relative Divergence of prominent media figures on the right using aligned word embeddings.

Salience of Clinton's e-mails and the Russia Investigation

In Figure 9, we notice that the term `email` first changes significance in period 8 and continues to change until period 11. The context stabilizes in periods 16-18. These changes track closely to those of `russia`, which is to be expected, given Russia's role in the controversy. The timing of these changes coincide with the investigations into Russian interference in the 2016 US elections. These terms both have a high maximum value of divergence, indicating the possibly polysemous nature in which they're being used. The more muted nature of `hack` is a stronger indicator of the e-mail controversy, owing to its usage in association with the John Podesta email hack. Compared to the use of `leak` which could refer to other administrative leaks with varying contexts in the Trump government, it has much greater variation in usage.

Figure 8 shows us how the context around the impeachment investigations have changed over time. We notice a spike in divergence in period 11 for `impeach`, right after the election, indicating how this was a hot-button issue immediately. Additionally, we notice a rise in divergence in period 13 for `hoax`, pointing to President Trump's claims calling Russian interference in the elections a `hoax` (noa, 2019). We see a higher and more disperse spread of divergence for `hunt`, beginning in period 17, when President Trump claimed that the Democrats were engaged in a 'witch hunt' (noa, 2018).

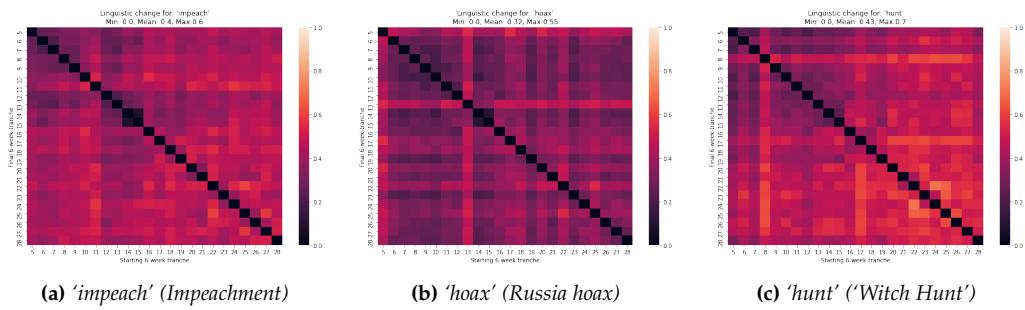


Figure 8: Relative Divergence of words related to Trump's scandals using aligned word embeddings.

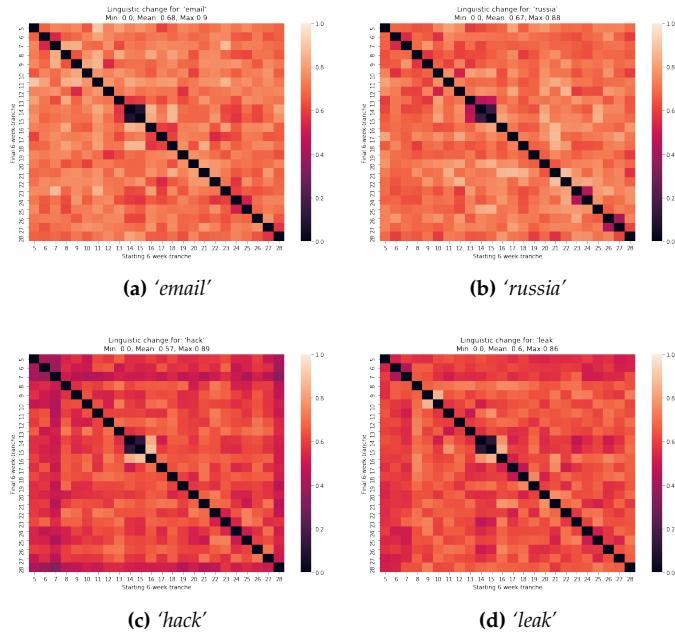


Figure 9: Relative Divergence of words related to Clinton’s e-mail controversy using aligned word embeddings.

Tracking tariffs and the trade war

President Trump’s focus on ‘Make America Great Again’ has also translated into renewed attention on American manufacturing, and the effect of free-trade agreements on jobs in America. As he tried to publicly push companies to move jobs back to the US, he also criticized past administrations for crafting unfavourable trade deals with China. This soon morphed into a trade war with China, with escalating tariffs levied by both sides. We notice the first change in the context surrounding the word *tariff* in period 13, in the first quarter of the presidency, followed by no change for an extended period. This is indicative of a significant change in rhetoric in period 13, until period 23. The change in period 23 is matched by those for words *trade* and *china*. This period included a serious escalation of the trade war with China.

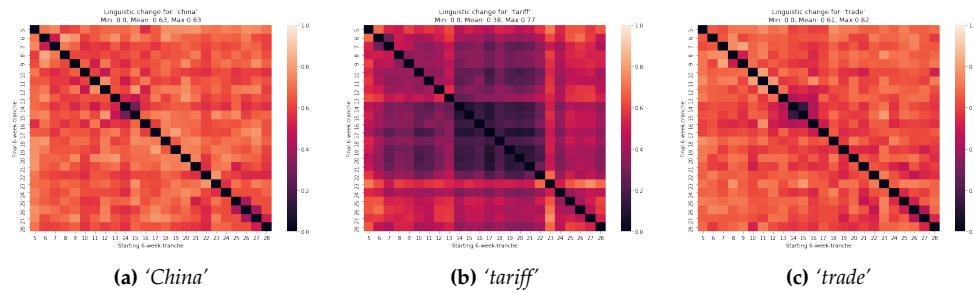


Figure 10: Relative Divergence of words related to Trump’s trade war with China using aligned word embeddings.

Using Projections of Word Embeddings to understand relative positions

Changing associations of political figures

Figure 11 provides a projection of political names and entities onto dimensions of security, trade, ideology, region, economy, and race. I select the names of the candidates in the 2016 US elections and the names of the political parties, to understand how their relative positions on these political dimensions change over time in the subreddit. This helps us understand how the community's language and context around these figures, and by extension their perception of these entities, has shifted over time.

On the 'Security' dimension in figure 11, we notice how `trump` ranks consistently high relative to all other entities. `clinton` moves towards the lower end of the dimension in the month of the election, indicating how the subreddit viewed her as being weak on Security issues. `democrat` is consistently on the lower dimension (weak security) in all but one window, that of the election. I interpret this as Clinton's fall from relevance as the opposition to President Trump after his election.

On questions of Ideology, we notice `trump` and even `democrat` not being fixed on the dimension, with only `republican` being static. We also notice how `bernie` and `stein` consistently represent the opposite of the Republican ideology. This latter observation should not be surprising – Bernie Sanders and Jill Stein have been purists for the duration of their political careers, and stand in stark opposition to the Republicans. Interestingly, in the period of the election, `stein` moves to the center while `clinton` moves to the extreme – a strong indicator of how perceptions in the lead-up to the elections were vastly different from other periods.

On questions of the Economy, `trump` loads heavily on the poor or non-elite end of the dimension, affirming the perception of his anti-establishment credentials. `bernie` occupies a similar position. `clinton` is positioned on the 'rich' end of the dimension in the month of the election, showcasing how her established pedigree came to the fore in that period. On the dimension of Race, `trump` shows up close to the extreme in two windows: Super Tuesday 2016, and August 2017 during the Charlottesville attack. Both these periods were dominated by strong racial and ethno-nationalist messaging by Trump and his staff.

Consistent perceptions of news outlets

Figure 12 provides a projection of news outlets – print and Television – on the aforementioned dimensions. I split them into these categories to disambiguate the different positions they take, and to examine equivalent pairs, like CNN-Fox, and NYT-Breitbart.

On the 'Security' dimension in figure 12a, Fox News consistently ranks high relative, with MSNBC and CNN occupying the 'weak' end of the dimension. The distance between Fox and the rest of the outlets from the second window onwards is indicative of the extreme posturing characteristic of their content. Fox also projects heavily on the 'conservative' end of the ideology dimension, but it's surprising to see CNN show up there in the third window (first quarter of the Presidency).

On matters of ‘Trade’, MSNBC and CNN load heavily on the ‘global’ end of the dimension. Other outlets like CBS and ABC promptly occupy the center on this dimension too. MSNBC and CNN occupy similar positions again on the ‘Economy’ dimension, with the community perceiving them to be more closely associated with the ‘poor’ end of the dimension. On questions of ‘Race’ again, CNN stays reliably to the ‘immigrant’ end of the dimension.

The consistency of the results persist with print news outlets too. National Review stakes its position as a conservative hawk on the Security dimension, its relative position far away from the other print outlets. The New York Times, Washington Post, and Wall Street Journal regularly load heavier on the ‘weak’ end of this dimension, with Breitbart closer to the Review. The conservative nature of the Review persists on the Trade dimension, along with more globalist positions taken by the Times, Post and Journal.

IV. DISCUSSION

The sharp increase in average cross-entropy for the community and top users in Period 10 is anomalous in several aspects. This period includes the weeks of the election, but it doesn’t have the highest number of posts. But it does have the highest number of unique posters in our time-frame, a sharp increase from the previous period, even if the successive period, with its lower cross-entropy, doesn’t see a sharp fall in posters (see Figure 1).

Breitbart’s position on the issue dimensions is an interesting case to examine. Given the consistency in the loadings of other outlets on these dimensions, one would expect Breitbart to have a relatively stable position across these windows. Instead, what we see is that Breitbart’s position shifts, closer at times to the conservative National Review and on the opposite end at others. One explanation for this could be that that constructed dimensions are not stable and don’t encode the information we think it does. Testing and tuning these dimensions for greater reliability would be essential in making these results generalizable. An alternate hypothesis is that the nature of Breitbart’s newsroom – digital, highly partisan, more finely attuned to and reliant upon online communities for amplification – make it fundamentally different from the more traditional outlets we compare against. Adding outlets like Vox Media, Huffington Post, or Buzzfeed could provide a third set of comparisons, among internet-first news outlets. But one could also examine how the language in Breitbart’s articles linked on the subreddit align with the language in the subreddit.

This paper examines how the linguistic characteristics of a politically-engaged community changes with time, with an examination of how key phrases shift in accordance to events. We’ve been able to examine the relationship between a user’s status and the language of the community, with our results replicating some of the findings in Danescu-Niculescu-Mizil et al.’s work. (Danescu-Niculescu-Mizil et al., 2013). Further work on this should focus on expanding beyond the n -gram model to quantify community language. Word Embedding methods which are able to better capture context could be used to identify more fine-tuned linguistic norms. Formalizing user status with centrality measures and sub-community network detection could help us tease out amplification channels within this community to examine in-group competition.



Figure 11: Projection of Candidate and Party names on dimensions of Security (strong-weak), Trade (local-global), Ideology (right-left), Region (urban-rural), Economy (rich-poor), and Race (american-immigrant).

From left to right, each projection refers to the time-window including (1) Super Tuesday in March 2016, (2) Elections in November 2016, (3) the first quarter of the presidency in Feb-Mar 2017, and (4) the 'Unite the Right' rally in Charlottesville, August 2017



(a) Television News Media. Includes 'MSNBC', 'CNN', 'ABC News', 'CBS News', and 'Fox News'.



(b) Print News Media. Includes 'New York Times', 'Washington Post', 'Wall Street Journal', 'National Review', and 'Breitbart'.

Figure 12: Projection of Television and Print media outlets on dimensions of Security (strong-weak), Trade (local-global), Ideology (right-left), Region (urban-rural), Economy (rich-poor), and Race (american-immigrant).
 From top to bottom, each projection refers to the time-window including (1) Super Tuesday in March 2016, (2) Elections in November 2016, (3) the first quarter of the presidency in Feb-Mar 2017, and (4) the 'Unite the Right' rally in Charlottesville, August 2017

REFERENCES

- Family separation under the Trump administration – a timeline | Southern Poverty Law Center. URL <https://www.splcenter.org/news/2019/09/24/family-separation-under-trump-administration-timeline>.
- Donald Trump Picks Reince Priebus as Chief of Staff and Stephen Bannon as Strategist. *The New York Times*, Nov. 2016. URL <https://www.nytimes.com/2016/11/14/us/politics/reince-priebus-chief-of-staff-donald-trump.html>.
- Timeline of the Muslim Ban, May 2017. URL <https://www.aclu-wa.org/pages/timeline-muslim-ban>. Library Catalog: www.aclu-wa.org.
- Donald Trump again rants against Robert Mueller 'witch hunt'. *USA Today*, 2018. URL <https://www.usatoday.com/story/news/politics/2018/04/10/donald-trump-rants-against-robert-mueller-witch-hunt/501973002/>.
- Trump and Putin talk of 'Russian hoax'. *BBC News*, May 2019. URL <https://www.bbc.com/news/world-us-canada-48141017>.
- D. C. Atkinson. Charlottesville and the alt-right: a turning point? *Politics, Groups, and Identities*, 6(2):309–315, Apr. 2018. ISSN 2156-5503. doi: 10.1080/21565503.2018.1454330. URL <https://doi.org/10.1080/21565503.2018.1454330>.
- J. Baumgartner, S. Zannettou, B. Keegan, M. Squire, and J. Blackburn. The Pushshift Reddit Dataset. *arXiv:2001.08435 [cs]*, Jan. 2020. URL <http://arxiv.org/abs/2001.08435>. arXiv: 2001.08435.
- B. Blodgett and A. Salter. Ghostbusters is For Boys: Understanding Geek Masculinity's Role in the Alt-right. *Commun Cult Crit*, 11(1):133–146, Mar. 2018. ISSN 1753-9129. doi: 10.1093/ccc/tcx003. URL <https://academic.oup.com/ccc/article/11/1/133/4953070>.
- C. Danescu-Niculescu-Mizil, R. West, D. Jurafsky, J. Leskovec, and C. Potts. No country for old members: user lifecycle and linguistic change in online communities. In *Proceedings of the 22nd international conference on World Wide Web - WWW '13*, pages 307–318, Rio de Janeiro, Brazil, 2013. ACM Press. ISBN 978-1-4503-2035-1. doi: 10.1145/2488388.2488416. URL <http://dl.acm.org/citation.cfm?doid=2488388.2488416>.
- J. Daniels. The Algorithmic Rise of the "Alt-Right". *Contexts*, 17(1):60–65, Feb. 2018. ISSN 1536-5042. doi: 10.1177/1536504218766547. URL <https://doi.org/10.1177/1536504218766547>.
- W. L. Hamilton, J. Leskovec, and D. Jurafsky. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. *arXiv:1605.09096 [cs]*, Oct. 2018. URL <http://arxiv.org/abs/1605.09096>. arXiv: 1605.09096.
- N. Heikkilä. Online Antagonism of the Alt-Right in the 2016 Election. *European journal of American studies*, 12(12-2), July 2017. ISSN 1991-9336. doi: 10.4000/ejas.12140. URL <http://journals.openedition.org/ejas/12140>.

-
- D. Jurafsky and J. H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall, 2 edition, 2008. ISBN 978-0-13-187321-6.
- T. Koulouris. Online misogyny and the alternative right: debating the undebatable. *Feminist Media Studies*, 18(4):750–761, July 2018. ISSN 1468-0777. doi: 10.1080/14680777.2018.1447428. URL <https://doi.org/10.1080/14680777.2018.1447428>.
- A. C. Kozlowski, M. Taddy, and J. A. Evans. The Geometry of Culture: Analyzing the Meanings of Class through Word Embeddings. *Am Sociol Rev*, 84(5):905–949, Oct. 2019. ISSN 0003-1224, 1939-8271. doi: 10.1177/0003122419877135. URL <http://journals.sagepub.com/doi/10.1177/0003122419877135>.
- A.-D. League. When Women are the Enemy: The Intersection of Misogyny and White Supremacy. Technical report, Center on Extremism, Anti-Defamation League, 2018. URL <https://www.adl.org/resources/reports/when-women-are-the-enemy-the-intersection-of-misogyny-and-white-supremacy>.
- A. Nagle. *Kill All Normies: Online Culture Wars From 4Chan and Tumblr to Trump and the Alt-Right*. Zero Books, 2017. ISBN 978-1-78535-544-8.
- D. Nieborg and M. Foxman. Mainstreaming Misogyny: The Beginning of the End and the End of the Beginning in Gamergate Coverage. In J. R. Vickery and T. Everbach, editors, *Mediating Misogyny: Gender, Technology, and Harassment*, pages 111–130. Springer International Publishing, Cham, 2018. ISBN 978-3-319-72917-6. URL https://doi.org/10.1007/978-3-319-72917-6_6. Type: 10.1007/978-3-319-72917-6_6.
- R. Řehůřek and P. Sojka. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010. ELRA. <http://is.muni.cz/publication/884893/en>.