

# Characterizing the Internet Research Agency’s Social Media Operations During the 2016 U.S. Presidential Election using Linguistic Analyses

Ryan L. Boyd<sup>1</sup>, Alexander Spangher<sup>2\*</sup>, Adam Fourney<sup>3</sup>, Besmira Nushi<sup>3</sup>,  
Gireeja Ranade<sup>3</sup>, James W. Pennebaker<sup>1</sup>, Eric Horvitz<sup>3</sup>

<sup>1</sup>University of Texas at Austin, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>Microsoft Research  
ryanboyd@utexas.edu, aas2230@columbia.edu, adamfo@microsoft.com, benushi@microsoft.com,  
gireeja@eecs.berkeley.edu, pennebaker@utexas.edu, horvitz@microsoft.com

## Abstract

Converging investigations on the part of multiple agencies/agents have provided overwhelming evidence for Russian interference in the 2016 U.S. presidential election. As a part (and consequence) of recent reports, multiple datasets that capture actions taken by actors of the Internet Research Agency (IRA), have been released to the public. In the current paper, we present an abridged report of several preliminary forensic analyses of Facebook ad data and Twitter troll accounts that were run by the IRA during the election cycle. Through the use of language analysis, we characterize the evolution of IRA content over the course of the election cycle, providing a basis for understanding how left- and right-leaning ideologies were differentially targeted to spread enmity among the American electorate. Additionally, through an analysis of syntactic constructions, we find that the content produced by the IRA on Twitter was linguistically unique from a control sample of English-speaking Twitter accounts. Altogether, our findings suggest that the IRA’s operations were largely unsophisticated and “low-budget” in nature, with no serious attempts at point-of-origin obfuscation being taken.

## Introduction

In May 2018, the Democratic representatives from the United States House Permanent Select Committee on Intelligence (USHPSCI) made public their findings regarding Russian interference in the 2016 United States presidential election. In their report, the USHPSCI supported and reaffirmed previous conclusions drawn by the Intelligence Community regarding widespread election interference taken by the Kremlin, ranging in scope from hacking-and-dumping campaigns to the dissemination of propaganda. Additionally, the Committee’s report revealed several details resulting from further investigation, including the pur-

chase and deployment of socially polarizing ads, webpages, and internet “trolls” by the Internet Research Agency (IRA), a Saint Petersburg-based company known to have engaged in long-term influence operations on behalf of Russian political and economic interests [1].

A consequence of the USHPSCI report has been the public release of two datasets reflecting the behaviors of IRA actors at the time of this writing. The first dataset includes over 3,500 Facebook advertisements purchased by the IRA – ads that were designed to fan the flames of discontent and spread general enmity within the American public [2]; it is estimated that over 11 million American internet users were exposed to these advertisements. The second dataset, which contains timelines for over 1,200 English-language Twitter accounts that were found to be operating as agents of the IRA, was curated and released in July 2018 by Darren Linvill and Patrick Warren [3]–[5]. Both datasets contain rich information in the form of timelines, behaviors, and the language used by IRA actors and associated metadata.

In this paper, we present an initial forensic analysis of the IRA data (i.e., Facebook ads and Twitter timelines) that has been publicly released at the time of this writing. Our primary goal is to simply characterize the social media behaviors of the IRA prior to, during, and following the 2016 election. It is our hope that these findings will serve to facilitate continued and deeper investigations of foreign interference in both past and future democratic processes. By demonstrating a small number of basic approaches to characterizing the IRA’s behaviors, we hope to assist the Intelligence Community, research community, and the general public in understanding the foundational details of when and how the IRA attempted to manipulate the psychological landscape surrounding the election.

---

\*Work performed during an internship at Microsoft Research, Redmond, WA.



# Analysis of IRA Facebook Advertisements

## Data and Methods Overview

All data was downloaded from the USHPSCI website, which provided ~3,500 PDF files containing the content of advertisements that were identified as having been purchased by the IRA [2]. All PDF files were subjected to OCR to extract the linguistic content and metadata pertaining to each ad (e.g., time of launch, number of click-throughs, amount of money spent per ad).

For this paper, we present two primary sets of analyses of the Facebook Ad data. The first is a simple forensic test to help confirm the origin of the ads, focusing on the time of day when ads were launched. Should the time of ad launches correspond to standard patterns of work behavior for the Saint Petersburg region, we will be able to determine whether any serious attempts at point-of-origin obfuscation occurred<sup>1</sup>.

The second analysis that we present is a general examination of the ad content over the course of the 2016 election cycle, allowing us to better understand how the content of the ads evolved over time. This analysis was conducted using the free open source software Meaning Extraction Helper [6] in conjunction with the statistical application *R* [7]. Put succinctly, we conducted a topic modeling procedure known as the meaning extraction method [8]–[11] to identify the most prominent themes across all of the Facebook Ad data, then quantified the degree to which each theme waxed and waned over the election cycle.

## Results

The time of each ad launch was converted to Moscow Standard Time (MST) and plotted as a distribution of number of ads launched, aggregated across the hours of each day. The ad launch versus time of day distribution is presented in Figure 1. Results from this analysis were quite striking, revealing that the overwhelming majority of ads were launched between 9:00am and 6:00pm MST – in essence, during standard business operation hours in Saint Petersburg. It does not appear that any attempts were made to obfuscate the point of origin by the use of delayed / automated launchers that would disperse the ads launch times across various latitudes.

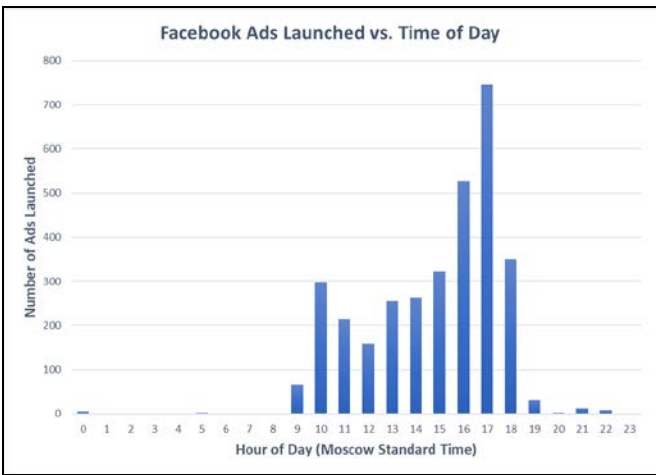


Figure 1

Number of Facebook Ads Launched per Hour of Day (Moscow Standard Time).

Note that, as part of the meaning extraction method, theme extraction and quantification were all performed automatically via principal components analysis; results were not manually selected from a larger pool of possible themes as is commonly done by researchers reporting results from methods like Latent Dirichlet Allocation<sup>2</sup>. For the purposes of the current research, we extracted the 10 themes from the Facebook Ad content that accounted for the largest percentage of variance in thematic content.

Theme labels and example words are presented in Table 1. Each ad was quantified for the degree to which it fit each of the 10 central themes using a regression approach [15]. Following quantification, each ad was classified with a binary score denoting that it was composed of either theme-related content (1) or no theme-related content (0). The binary score was assigned using a percentile-based cut-off (1.5 standard deviations above each theme’s mean, corresponding to roughly the 85<sup>th</sup> percentile). This threshold was selected to reduce statistical noise to above-ambient levels. Changes to the selected threshold in either (i.e., lower versus higher thresholds) had no substantive impact on the results presented in this paper. Following the quantification / assignment of each ad to the 10 most prominent ad themes, the distribution of ad launch dates were plotted over time, separately by theme<sup>3</sup>.

<sup>1</sup> Although not discussed further in this paper, we also note that a substantial percentage of the ads were paid for using Russian currency (i.e., rubles), which serves as an additional hint that systematic measures to obfuscate the point-of-origin may not have been taken by the IRA. Future analyses may consider determining if obfuscation techniques were employed inconsistently (e.g., time of ad purchase covarying with currency, etc.).

<sup>2</sup> While many variations of topic modeling procedures exist, we employ the meaning extraction method due to its use of principal components analysis. This allows us to hone in on a small number of themes that are responsible for the largest amount of thematic variance within the dataset [12]–[14], resulting in a small number of meaningful topics, resulting in our ability to synthesize an interpretation that would not be possible with hundreds or thousands of topics.

<sup>3</sup> Smoothed density plots generated using the “ggridges” package in R [16].

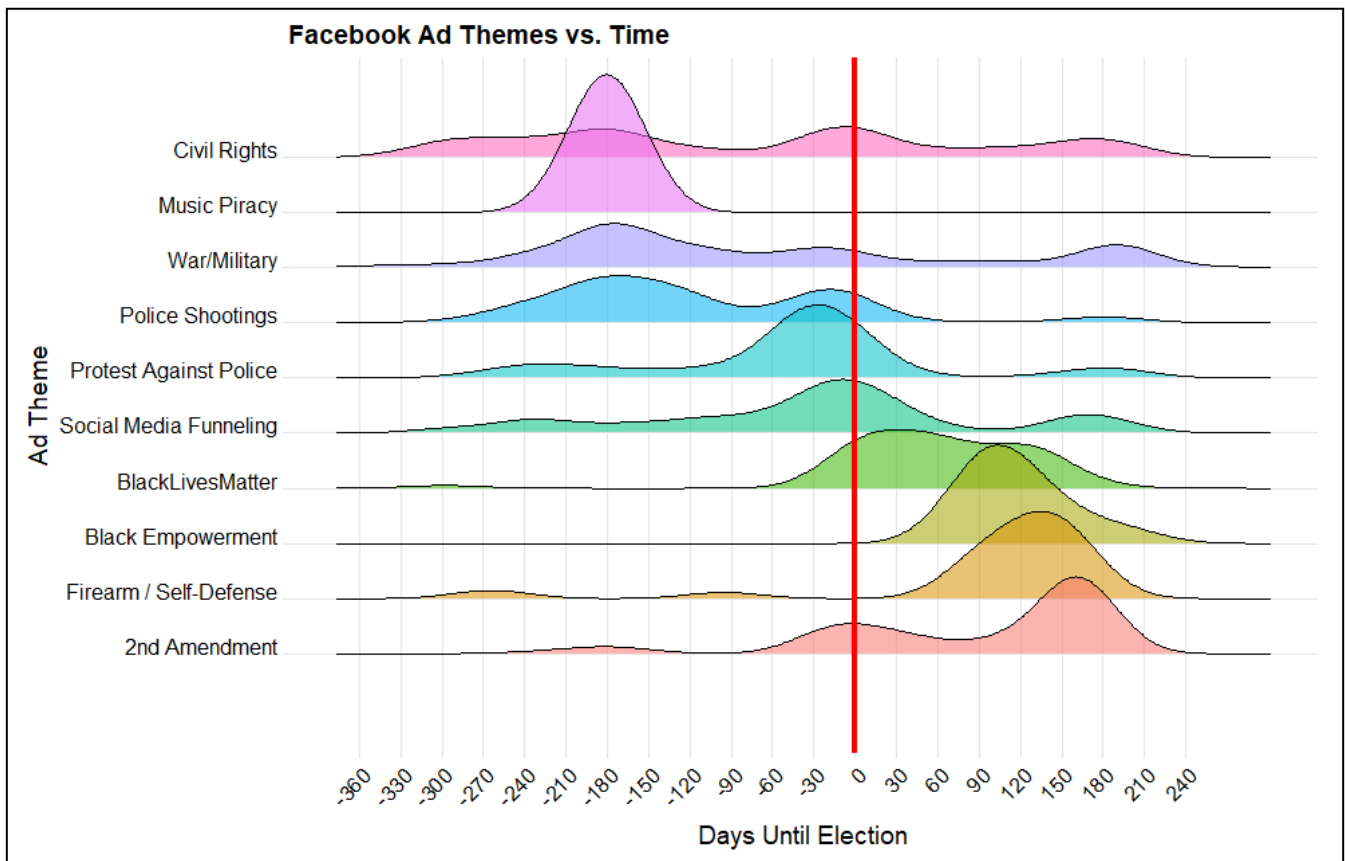


Figure 2

*Evolution of Content in IRA Facebook Ads Before, During, and After the 2016 U.S. Presidential Election.*

*Notes: The vertical red line denotes election day (November 8<sup>th</sup>, 2016). The Y axis reflects the number of ads launched each day that are composed to some degree of the respective themes; each tick mark corresponds to roughly 30 ads.*

Results from the above-described procedure are presented in Figure 2. We draw the readers’ attention to some key features of the thematic evolution found in the IRA Facebook Ad campaign. First, and as reported elsewhere [17], the majority of thematic content was engineered to seed discord among individuals in the U.S. with differing ideologies. However, and notably, other tactics were employed as well, such as the deployment of ads loading highly on the Social Media Funneling and Music Piracy themes, designed to drive off-site page clicks to other IRA-sponsored content (i.e., websites and IRA accounts not hosted on Facebook).

In looking at how the themes unfolded over time, one can see that the thematic content appeared in two forms: persistent and provisional. As an example, the *Civil Rights* theme appeared well before the election and persisted to varying degrees until the end of the timespan covered by the data.

Theme Label	Example Words
Civil Rights	Equality, civil, right
Music Piracy	Music, free, download, online
War/Military	Syria, terrorist, war, military
Police Shootings	Officer, cop, shoot, fatally
Protest Against Police	October 22, protest, brutality
Social Media Funneling	Follow, subscribe, channel
BlackLivesMatter	#blacklivesmatter, #unite4justice
Black Empowerment	#blackpower, #africanandproud
Firearm / Self-defense	Self-defense, event, safe, NY
2 <sup>nd</sup> Amendment	2 <sup>nd</sup> , amendment, patriot, gun

Table 1

*Top 10 Themes Extracted from the Facebook Ad Data. Theme labels were manually applied based on observing ads that loaded highly on each theme. A complete component loading list is publicly available at <https://osf.io/zcyt7>*

In contrast, themes such as *Protest Against Police* and *BlackLivesMatter* appeared quite rapidly, drastically scaling up shortly before election day and disappearing after the conclusion of the election.

A second notable feature of how the ads' thematic content evolved is that, at any given time point, efforts appeared to primarily target one end of the political spectrum rather than both simultaneously. For example, the focus on *War/Military* is largely concerned with a wide range of military issues, as is illustrated in the following excerpt from one of the highest-loading ads for this theme:

*the u.s. flag hasn't been lowered in honor of the murdered marines [...] by doing so obama shows the highest degree of disrespect for the fallen soldiers. it turns out that for him the soldiers as well as veterans do not deserve any attention or respect.*

Following the decline of the *Military/War* theme, we see a long-term focus on issues primarily attended to by the ideological left, including an explicit distrust/disdain for the police and high fixation on issues related to Black identity and the “Black Lives Matter” movement. The above excerpt can be contrasted with text sampled from one of the top-loading ads for the *BlackLivesMatter* theme, designed to exploit racial tensions in the U.S.:

*overheard two people (white old ladies) talking about a particular film they watched recently. their comments about black female actress almost made me choke on my saliva. believe me, i had to step-in, cus these ugly madafakas were describing our queens as if it was the fault of black female actors to be dope and on-point af.*

Following the election, the focus shifted back again to primarily conservative issues, such as self-defense and gun rights. We note, however, at most points within +/-100 days of the election, some combination of themes targeting both the ideological left and the ideological right were present.

## Analysis of IRA Twitter Troll Accounts

### Data and Methods Overview

IRA troll account data was downloaded from the FiftyEight GitHub repository on August 21, 2018 [18]<sup>4</sup>. As with the Facebook Ad data, we were first interested in de-

<sup>4</sup> Since our original download of the data, updates have been pushed to the GitHub master repository. These updates include some additional cleaning of the data, such as removing duplicate entries. None of the updates to the source dataset affect the results / conclusions of this project.

termining the temporal patterns of troll activities and exploring whether account behaviors mapped onto what might be expected during standard business hours in Saint Petersburg. Rather than look directly at the number of tweets posted during each hour of the day across all accounts, we instead summarize the data at the *account* level (rather than the individual *tweet* level). This was done to prevent undue influence from high-volume accounts. As such, the number of tweets made during each hour of the day was normalized in a within-subject fashion prior to aggregation.

In contrast to the text analysis methods used with the Facebook Ad data, we turned to a more intensive natural language processing/machine learning joint methodology to create a robust forensic test of the tweets' origins. The primary question that we seek to address with this analysis is: were the IRA troll accounts creating unique content or, alternatively, simply parroting native, non-troll tweets? Put another way, to what extent did the IRA troll accounts *seed* discontent versus simply *amplify*<sup>5</sup> discontent that was already in circulation?

The methods used for this forensic approach were significantly more advanced than those performed on the Facebook Ad data. Briefly described, the analyses performed were designed to robustly test whether the linguistic signature of the IRA Twitter accounts was congruent with organic, native English content or, instead, syntactically distinct from tweets made by the general English-speaking West. To conduct this test, we compared the syntactic patterns of the IRA troll accounts to a random sampling of archival, English-speaking Twitter accounts that were warehoused in the first author's collection; the control sample was selected to include tweets made during the same period as the IRA's troll account activity.

Rather than simply testing the degree of differentiability in the language of the IRA sample versus the Control sample, such an analysis must “work backwards” through the language of each sample to determine differentiability. The fact that a random forest algorithm, for example, can separate the linguistic signature of the IRA accounts from the troll accounts is not particularly meaningful<sup>6</sup>. Such separability could be driven by several factors, such as superficial stylistics (e.g., punctuated, “news headline” style language versus casual social media sharing). However, as we iteratively strip away such superficial differences, user accounts from each sample should converge into an inseparable pool were they to have come from the same general population. This is particularly true to syntactic constructions (i.e.,

<sup>5</sup> Importantly, we are not interested in retweets for these linguistic analyses, as these constitute explicit amplification behaviors. Instead, we investigate here the linguistic patterns of IRA tweets that are not explicit amplification of signals. As noted in our discussion, such an analysis is better suited to an information spread / social network style of approach.

<sup>6</sup> 10-fold cross-validated random forest Cohen's  $\kappa = 0.89$ .



syntactic n-grams, or sn-grams, as opposed to lexicon-driven n-grams), and has become a widely-adopted design in attribution methodologies [19].

Put in simple terms, as we try to intentionally force convergence between the linguistic patterns found in the IRA accounts and Control accounts, respectively, we will see one of two outcomes. Were the IRA accounts amplifying the visibility of native English content, the two samples should rapidly converge into a single pool that cannot be differentiated. On the other hand, were the IRA accounts generating unique content on their own, the two samples should remain separable despite iterative attempts to force a lack of differentiability. For these analyses, we adapted and innovated on established methods described elsewhere referred to as the “unmasking” method [20]–[23].

For these analyses, we considered only IRA accounts that posted tweets in the English language. We relied on the metadata accompanying each tweet in the public dataset to filter out non-English tweets; the same procedure was used to filter the control sample. Standard cleaning procedures were applied via regular expressions, including the removal of URLs<sup>7</sup>, retweets, and hashtags, and the replacement of usernames with proper nouns. A minimum criterion for inclusion was set at 100 tweets (temporal analysis) or 100 total tokens (linguistic analyses). All tweets were aggregated by account, resulting in a troll account sample size of  $N=969$  accounts ( $M$  token count = 2157.10;  $SD$  = 5846.48; min = 102; max = 107484) and final control sample size of  $N=1078$  accounts ( $M$  token count = 5245.65;  $SD$  = 9433.31; min = 101; max = 78432). All measures (e.g., percent of each text comprised of sn-grams) were normalized by each account’s total number of tokens, effectively controlling for differing activity levels between accounts.

Following text cleaning, all content was tagged for Part of Speech (POS) using Stanford’s CoreNLP framework [24], running the GATE Twitter POS model [25]. Following POS tagging, we extracted s1-grams through s3-grams for all accounts, retaining the top ~1000 sn-grams that were common across both samples (top 1000 determined independently for each sample, then matched across samples).

## Results

As with the Facebook Ad data, we present the average IRA user activity distribution (converted to Moscow Standard Time) in Figure 3. Parallel to the findings reported earlier, we see again that the majority of IRA troll activity occurred between the hours of 9:00am and 6:00pm MST. Additionally, user activity peaked at 5:00pm – identical to the peak seen in the Facebook Ad data. Here again we see

<sup>7</sup> All tweets containing URLs were omitted to prevent contamination from text generated as part of link previews. Inclusion of these tweets did not alter the results.

no serious attempt to obfuscate the activity patterns of IRA actors. However, the patterns for troll account activity was noticeably noisier than that seen with the Facebook ads, and we do witness account activity across all hours of the day.

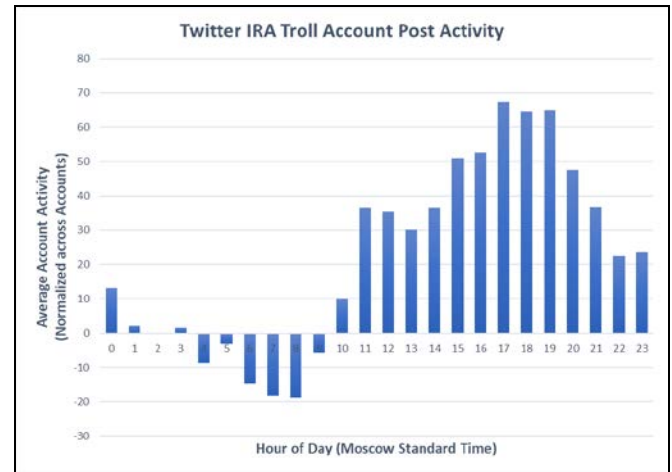


Figure 3

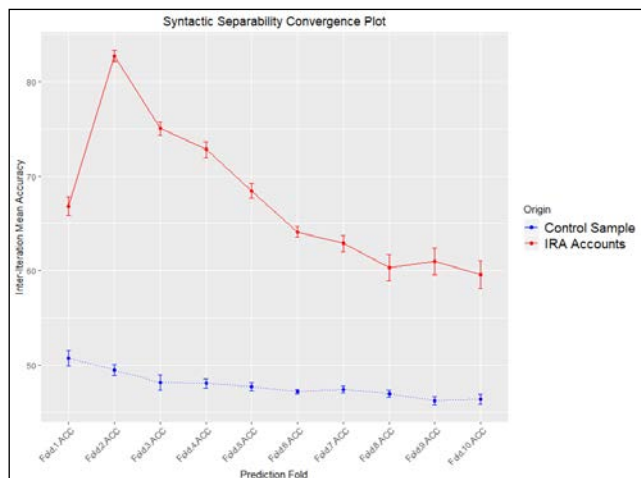
*Normalized, Relative IRA Twitter Account Activity per Hour of Day (Moscow Standard Time).*

Turning now to our analysis of syntactic patterns, the analysis of IRA account syntax revealed a strong differentiation trend that persisted despite attempts to suppress linguistic uniqueness. Results are presented in Figure 4. As a brief interpretation of the results, first consider the bottom, dotted line, representing the accuracy achieved when attempting to distinguish random subsets of the control sample from itself. Throughout the entire analysis, an accuracy of roughly 50% persisted across all folds, and all folds across all iterations – no better than chance<sup>8</sup>. On the other hand, the solid, top line represents the accuracy achieved when attempting to distinguish the IRA accounts from random subsamples of the control group. Were the two samples to have come from the same broader population of accounts, we would have expected both to rapidly converge to 50% accuracy (or lower).

In short, the syntactic constructions of the IRA troll account language fail to converge with general English-language Twitter accounts, suggesting a distinct regional origin where English is a non-native language. Further analyses could be performed to identify the specific differentiating syntactic patterns to better confirm syntactic transfer from Russian L1 speakers to English as an L2 language.

<sup>8</sup> As part of the methods employed here for these analyses, class-balancing is performed prior to differentiation. As a result, the 50% accuracy baseline corresponds to intuition on probabilities. We refer readers to the references listed earlier for the “unmasking” method for a more detailed description of procedures employed in this analysis.

Figure 4



Syntactic Separability plot. Results are averaged across 10 folds, which were themselves averaged across 100 iterations. The X-axis represents each of the 10 folds; the Y-axis represents the average accuracy score across all iterations.

## Post Hoc Descriptive Analyses

Beyond the “unmasking” approach just reported, we conducted several additional, but preliminary, follow-up analyses to better understand the strongest differentiating factors among the language patterns found in the IRA and control Twitter account samples. These *post hoc* contrasts were conducted with the goal of better understanding differences in both linguistic and psychological content. For these follow-up analyses, we conducted a robust pairing of inferential statistics with machine learning practices. Unlike the unmasking methods performed above, which worked “backwards” against differentiability, the *post hoc* analyses reported here work “forwards” to find the most discriminating features among the two samples.

We urge the reader to treat all *post hoc* conclusions tentatively in lieu of additional, more rigorous study.

### Bootstrap Aggregated *t*-tests

Our initial analysis to highlight the most distinguishing features was performed as a series of independent, bagged Student’s *t*-tests. We specifically tested both the sn-grams that were detailed in the previous section, as well as features extracted using the LIWC2015 dictionary [26]. The results of these models were also subjected to the highly

Feature Origin	Feature	Mean <i>t</i> Value	Mean <i>p</i> Value, Bonferroni Adjusted	Control Sample <i>M</i>	IRA Sample <i>M</i>	Example POS-tagged n-gram (%DIFF)
sn-gram	UH PRP	9.54	< 0.01	0.59	0.06	🤖/UH i/PRP (9611.90)
sn-gram	UH	9.26	< 0.01	3.98	0.77	🤖/UH (13039.63)
LIWC2015	Authentic	9.37	< 0.01	58.93	29.70	
sn-gram	NN UH	8.64	< 0.01	0.56	0.06	time/NN lol/UH (879.35)
LIWC2015	informal	8.43	< 0.01	4.39	1.70	
LIWC2015	Sixltr	-8.91	< 0.01	13.43	19.07	
LIWC2015	power	-8.68	< 0.01	2.34	4.10	
LIWC2015	netspeak	7.65	< 0.01	2.40	0.79	
sn-gram	UH UH	6.94	< 0.01	1.48	0.14	😄/UH 😄/UH (9652.71)
LIWC2015	i	8.00	< 0.01	6.32	3.02	
LIWC2015	time	7.20	< 0.01	6.08	4.41	
sn-gram	NNS	-7.11	< 0.01	3.19	4.54	reforms/NNS (-96.45)
sn-gram	RB VBD	7.20	< 0.01	0.30	0.12	just/RB ate/VBD (2185.15)
sn-gram	: NNP	-6.95	< 0.01	0.10	0.53	:/: Trump/NNP (-86.14)
LIWC2015	Clout	-7.06	< 0.01	52.75	69.87	
sn-gram	PRP RB	6.86	< 0.01	0.70	0.34	i/PRP kinda/RB (3246.12)
sn-gram	. UH	6.88	< 0.01	0.74	0.20	./ 😄/UH (3327.73)
sn-gram	PRP	6.17	0.01	8.74	6.18	he/PRP (5.49)
LIWC2015	we	-6.25	0.01	0.65	1.39	
sn-gram	RB	5.93	0.01	5.91	4.53	lowkey/RB (3545.36)
LIWC2015	adverb	5.97	0.02	5.60	4.21	
LIWC2015	risk	-5.95	0.02	0.53	0.99	
LIWC2015	focuspast	6.35	0.02	3.12	2.14	
LIWC2015	ppron	5.89	0.02	11.30	8.14	
LIWC2015	swear	5.66	0.03	1.30	0.44	

Table 2

Top 25 Most Distinguishing Features (sn-grams and LIWC2015) Among the Twitter IRA and Control Samples All values are aggregated across bootstrapping iterations. Descriptions and examples of Part of Speech tags can be found at <https://www.clips.uantwerpen.be/pages/mbp-tags>



conservative Bonferroni adjustment to further reduce  $p$ -value inflation and control Type I error rates [27]. As such,  $p$ -values for these analyses should be considered as *extremely* conservative probability estimates.

Table 2 contains results for the top 25 most distinguishing features resulting from the bagged  $t$ -test analyses. We remind readers that means for each group are presented at the *account* level rather than at the tweet level. Because of the often non-intuitive nature of sn-grams, we also present in Table 2 examples of differentiating POS-tagged n-grams along with their %DIFF keyness scores [28], to assist in the interpretation of these results<sup>9</sup>.

We find that the strongest differentiating sn-grams often were highly stylistic, such as the relatively high use “utterances” (primarily emoji, plus words like “wow” and “yeah”) in the control sample versus almost none in the IRA sample. Similarly, use of personal pronouns, both in general and as part of several syntactic constructions, were typically used in the IRA sample at roughly half the rate as found in the control sample. Interestingly, we do not find differentiation among the two groups in their rates of articles or determiners (e.g., the, a, an; all  $ps > 0.99$ ), which are commonly cited as a highly visible markers of Russian-native English-learned syntactic transfer. However, our analyses did not explore possible mis-selection of determiners based on definiteness (e.g., using the word “the” instead of “a”), which are often seen across levels of fluency throughout second language learning [29].

Regarding the LIWC2015 features, it is interesting that one of the two summary measures [26] that were identified as the largest difference among the two samples was the *Authentic* measure, an index of deception [30]. Briefly, the *Authentic* measure is a psychological metric that captures the degree to which language is spontaneous and unfiltered, as measured through established patterns of language use (high scores = more spontaneous, low scores = more cautious and constructed). The IRA troll accounts exhibited extremely low *Authenticity* relative to the control sample. The degree to which these preliminary results are driven by the unique syntactic transfers present in the IRA accounts are currently unknown and, as such, should be interpreted with caution.

Altogether, the results from these linguistic analyses provide evidence that the IRA actors were not only composing their own tweets, but were doing so in a carefully constructed, intentionally deceptive manner. Importantly, the stylistic composition of the IRA tweets were unique and consistently differentiable from a general population sample in the same language.

## Discussion

---

<sup>9</sup> Keyness scores were calculated as a function of the percent of each sample that used each POS-tagged n-gram at least once rather than raw frequency scores in accordance with the notion of differentiating *individuals* rather than corpora more broadly.

In this paper, we presented initial analyses of behavioral patterns, thematic content, and syntactic/psychological constructions of Facebook ads and Twitter content propagated by the Internet Research Agency on behalf of Russian political interests. Our analyses provide confirmatory evidence of Russian social media behaviors intended to influence the 2016 U.S. presidential election and, additionally, provide a basis for characterizing the nature of such behaviors insofar as they were designed to alter the psychological element of the U.S. electorate.

One of the key findings of this paper is the fact that the IRA does not appear to have undertaken any serious attempts at concealing the point-of-origin for their Facebook ads or Twitter accounts. For both datasets, the activity patterns map squarely onto standard business hours in Moscow Standard Time. Moreover, the linguistic patterns of the Twitter activity show high differentiability from organic English-language accounts, including a high degree of deceptive psycholinguistic patterning. These findings suggest that the Russian/IRA approach was likely a low-budget, fairly blunt approach to disinformation / influence operations. Rather than employing complex or intricate techniques to cover their tracks, the IRA appears to instead have relied on a broad coverage approach, attempting to spread a high amount of polarizing content while minimizing expenditure.

The second key finding is more descriptive. An analysis of thematic content found in the IRA’s Facebook ad campaign allows us to better understand the timeline of targeting that occurred during the election cycle. The particular focus on disseminating ads designed to frustrate the ideological left in close temporal proximity to election day may be suggestive of the intent to disenfranchise (and thus suppress) left-leaning voters; further research must be conducted to explore this idea.

Critically, there is much room for future research/investigation into the matter of interference in the 2016 presidential election. Recent work on frame identification, for example, could help to more deeply understand how the Facebook ads and troll accounts may have developed their targeting strategies as a function of both U.S. and Russian media [31]. Deeper analyses of other types of metadata (e.g., number of click-throughs on Facebook ads, social network / retweeting patterns in the Twitter data) will help us to understand the extent to which the influence operations were effective in polarizing and/or suppressing the American electorate.

## Conclusion

As members of an increasingly global society, it is critical that researchers across disciplines investigate threats to the common good. With continued research, we will be able to better anticipate and identify future meddling in the democratic process, as well as international influence operations



more generally. Investigators around the world are continuing to make discoveries about the nature of the 2016 U.S. presidential election interference, creating opportunities for additional study, reflection, and planning. We call on members of the international research and intelligence communities to assist in unmasking those who attempt to undermine or harm the basic human rights of self-determination, freedom, and fairness.

## Acknowledgment

The authors acknowledge all individuals whose investigations have helped to discover and identify sources of interference in the U.S. election process, including (but by no means limited to) Special Counsel Robert S. Mueller III, members of the U.S. Department of Homeland Security, Darren Linvill, Patrick Warren, the international Intelligence Community, and the United States House Permanent Select Committee on Intelligence.

Preparation of this manuscript was aided by grants from the National Institute of Health 5R01GM112697-02), John Templeton Foundation (#48503) and the National Science Foundation (IIS-1344257). The views, opinions, and findings contained in this paper are those of the author(s) and should not be construed as position, policy, or decision of the aforementioned agencies, unless so designated by other documents.

## References

- [1] USHPSCI, "Exposing Russia's Effort to Sow Discord Online: The Internet Research Agency and Advertisements," *Exposing Russia's Effort to Sow Discord Online: The Internet Research Agency and Advertisements | U.S. House of Representatives*, 2018. [Online]. Available: <https://democrats-intelligence.house.gov/social-media-content/>. [Accessed: 31-Aug-2018].
- [2] USHPSCI, "Social Media Advertisements," 2018. [Online]. Available: <https://democrats-intelligence.house.gov/social-media-content/social-media-advertisements.htm>. [Accessed: 31-Aug-2018].
- [3] O. Roeder, "Why We're Sharing 3 Million Russian Troll Tweets," *FiveThirtyEight*, 31-Jul-2018. .
- [4] M. Staton, "Faculty measure impact of underreported activity of political Twitter trolls," *Newsstand | Clemson University News and Stories, South Carolina*, 2018. [Online]. Available: <http://newsstand.clemson.edu/faculty-measure-impact-of-underreported-activity-of-political-twitter-trolls/>. [Accessed: 21-Aug-2018].
- [5] D. L. Linvill and P. L. Warren, "Troll Factories: The Internet Research Agency and State-Sponsored Agenda Building," 2018.
- [6] R. L. Boyd, *MEH: Meaning Extraction Helper [Software]*. 2018.
- [7] R Core Team, *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2018.
- [8] R. L. Boyd, S. R. Wilson, J. W. Pennebaker, M. Kosinski, D. J. Stillwell, and R. Mihalcea, "Values in Words: Using Language to Evaluate and Understand Personal Values," in *Proceedings of the Ninth International AAAI Conference on Web and Social Media*, 2015, pp. 31–40.
- [9] C. K. Chung and J. W. Pennebaker, "Revealing Dimensions of Thinking in Open-Ended Self-Descriptions: An Automated Meaning Extraction Method for Natural Language," *Journal of Research in Personality*, vol. 42, no. 1, pp. 96–132, Feb. 2008.
- [10] A. Kramer and C. Chung, "Dimensions of Self-Expression in Facebook Status Updates," in *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, 2011, pp. 169–176.
- [11] N. Ramirez-esparza, C. K. Chung, E. Kacewicz, and J. W. Pennebaker, "The psychology of word use in depression forums in English and in Spanish: Testing two text analytic approaches," in *In Proc. ICWSM 2008*, 2008.
- [12] I. Jolliffe, "Principal Component Analysis," in *International Encyclopedia of Statistical Science*, M. Lovric, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 1094–1096.
- [13] J. Shlens, "A Tutorial on Principal Component Analysis," *arXiv:1404.1100 [cs, stat]*, Apr. 2014.
- [14] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics and Intelligent Laboratory Systems*, vol. 2, no. 1, pp. 37–52, Aug. 1987.
- [15] R. L. Boyd, "Psychological text analysis in the digital humanities," in *Data analytics in the digital humanities*, S. Hai-Jew, Ed. New York: Springer International Publishing, 2017, pp. 161–189.
- [16] C. Wilke, *Geoms to make ridgeline plots with ggplot2. Contribute to clauswilke/gggridges development by creating an account on GitHub*. 2018.
- [17] A. Ng, "This was the most viewed Facebook ad bought by Russian trolls," *CNET*, 2018. [Online]. Available: <https://www.cnet.com/news/this-was-the-most-viewed-facebook-ad-bought-by-russian-trolls/>. [Accessed: 03-Sep-2018].
- [18] FiveThirtyEight, "3 million Russian troll tweets," 2018. [Online]. Available: <https://github.com/fivethirtyeight/russian-troll-tweets>. [Accessed: 21-Aug-2018].
- [19] G. Sidorov, F. Velasquez, E. Stamatatos, A. Gelbukh, and L. Chanona-Hernández, "Syntactic N-grams as machine learning features for natural language processing," *Expert Systems with Applications*, vol. 41, no. 3, pp. 853–860, Feb. 2014.
- [20] R. L. Boyd, "Mental profile mapping: A psychological single-candidate authorship attribution method," *PLOS ONE*, vol. 13, no. 7, p. e0200588, Jul. 2018.
- [21] M. Coulthard, A. Johnson, D. Wright, A. Johnson, and D. Wright, *An Introduction to Forensic Linguistics : Language in Evidence*. Routledge, 2016.
- [22] M. Koppel, J. Schler, and E. Bonchek-Dokow, "Measuring differentiability: Unmasking pseudonymous authors," *Journal of Machine Learning Research*, vol. 8, no. 2007, pp. 1261–1276, 2007.
- [23] M. Koppel and Y. Winter, "Determining if two documents are written by the same author," *Journal of the American Society for Information Science and Technology*, vol. 65, no. 1, pp. 178–187, 2014.





- [24] C. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. Bethard, and D. McClosky, *The Stanford CoreNLP Natural Language Processing Toolkit*. 2014.
- [25] L. Derczynski, A. Ritter, S. Clark, and K. Bontcheva, “Twitter Part-of-Speech Tagging for All: Overcoming Sparse and Noisy Data,” in *Proceedings of the International Conference on Recent Advances in Natural Language Processing*, 2013.
- [26] J. W. Pennebaker, R. L. Boyd, K. Jordan, and K. Blackburn, “The development and psychometric properties of LIWC2015,” Austin, TX, 2015.
- [27] S. R. Narum, “Beyond Bonferroni: Less conservative analyses for conservation genetics,” *Conserv Genet*, vol. 7, no. 5, pp. 783–787, Oct. 2006.
- [28] C. Gabrielatos and A. Marchi, “Keyness: Appropriate metrics and practical issues,” in *Proceedings of the 2012 International Conference on Corpus-assisted Discourse Studies*, University of Bologna, Italy, 2012.
- [29] T. Ionin, M. L. Zubizarreta, and S. B. Maldonado, “Sources of linguistic knowledge in the second language acquisition of English articles,” *Lingua*, vol. 118, no. 4, pp. 554–576, 2008.
- [30] M. L. Newman, J. W. Pennebaker, D. S. Berry, and J. M. Richards, “Lying words: Predicting deception from linguistic styles,” *Personality and Social Psychology Bulletin*, vol. 29, no. 5, pp. 665–675, 2003.
- [31] A. Field, D. Kliger, S. Wintner, J. Pan, D. Jurafsky, and Y. Tsvetkov, “Framing and Agenda-setting in Russian News: a Computational Analysis of Intricate Political Strategies,” *arXiv:1808.09386 [cs]*, Aug. 2018.