# From Facebook to the Streets: Russian Troll Ads and Black Lives Matter Protests

Ugochukwu Etudo
School of Business
University of Connecticut
ugochukwu.etudo@uconn.edu

Victoria Yoon
School of Business
Virginia Commonwealth
University
vyyoon@vcu.edu

Niam Yaraghi
School of Business
University of Connecticut
niam.yaraghi@uconn.edu

## Abstract

*Online trolling is typically studied in the IS literature as an uncoordinated, anarchic activity. Coordinated, strategic online trolling is not well understood despite its prevalence on social media. To shed light on this prevailing activity, the present study examines the proposition that coordinated online trolling is timed to leverage macro societal unrest. In testing this proposition, we analyzed the dynamics of the Russian State's coordinated trolling campaign against the United States beginning in 2015. Using the May 2018 release of all Russian Troll Facebook advertisements, this study constructs a topic model of the content of these ads. The relationship between ad topics and the frequency of Black Lives Matter protests is examined. We argue that the frequency of Black Lives Matter protests proxies for civil unrest and divisiveness in the United States. The study finds that Russian ads related to police brutality were issued to coincide with periods of higher unrest. This work also finds that during periods of relative calm (evidenced by lower frequency of protests) Russian ads were relatively innocuous.*

## 1. Introduction

Online trolling is prevalent on social platforms. While academicians have produced scattered definitions of online trolling, the provenance of the term is such that the articulation of a stable and generally accepted definition may remain elusive. One definition states that online trolling is "the practice of behaving in a deceptive, destructive or disruptive manner in a social setting on the internet with no apparent instrumental purpose" [2, p. 97]. This definition implies that trolling is inherently anarchic. However, the trolling referenced in this study were well coordinated and had very specific, non-anarchic goals.

In a January 6th 2017 Intelligence Community Assessment (ICA) [1], it was announced that the United States intelligence community, in a rare and unanimous opinion, believed that the Russian government attempted to interfere with the United States presidential election and electoral process. The report, ICA 2017-01D, is a coordinated assessment involving intelligence information collected by The Central Intelligence Agency (CIA), The Federal Bureau of Investigation (FBI) and The National Security Agency (NSA). It is important to note that the report does not cover the effects of any purported activity by Russian state and state-sponsored actors, but describes that activity. At the time of the assessment, the intelligence community had sufficient evidence to report that "Russia used trolls as… part of its influence efforts to denigrate Secretary Clinton. This effort amplified stories on scandals about Secretary Clinton and the role of WikiLeaks in the election campaign" [1, p. 4]. According the ICA, a St. Petersburg based entity, "Internet Research Agency," (IRA) described in the report as consisting primarily of professional trolls spear-headed this effort.

Since the publication of ICA 2017-01D, additional details regarding Russian state and state-sponsored activities to subvert ordinary democratic processes in the United States have been released with increasing specificity. The U.S. Congress, in the final months of 2017, released a list of 2,752 (since deactivated and scrubbed) Twitter handles that Twitter Inc. flagged as emanating from Russia's Internet Research Agency. Shortly after, in February of 2018, the textual content of thousands of Tweets issued by these accounts is released. In November of 2017 the House Select Committee on Intelligence provided the public with a sample of Facebook advertisements created by the IRA. Finally, in May of 2018 congress published the comprehensive list of Facebook advertisements (3,393) and metadata pertaining to them. An often-cited objective of these trolling campaigns is to sow

HICSS

discord in the voting populace by amplifying voices on different sides of polarizing issues. This is the "trolling" activity referenced in the ICA. In particular (and we show this later) the Russian IRA, when creating Facebook advertisements, homed in on the Black Lives Matter (BLM) movement. The ICA refers to a subset of tactics used by Russia as *trolling*. The "trolling" activity described in the ICA is a coordinated activity, planned outside of the targeted platform, and executed in ordered and strategic fashion. However, much of the IS literature is focused on trolling as individual, uncoordinated behaviors for the amusement of the individual perpetrators [3]. Coordinated, strategic online trolling is not well understood. This research addresses the gap in academic coverage of online trolling where a central actor coordinates the activity of multiple virtual personas and entities.

This study aims to shed light on coordinated and non-anarchic online trolling. To do so, we examine the commonly held view that Russian troll activity is (Russian trolls continue to be very active) intended to sow discord within the American voting public. We analyze over 3,000 individual PDF documents published by the Minority in the House Intelligence Committee where each document corresponds to a single IRA ad buy. We describe the challenges in analyzing the data given its published format. We combine this data into a time series with the only publicly available collection of Black Lives Matter movement related protests. Using this novel dataset, we examine the strategy of IRA ad buys, testing the conventional wisdom that these buys were intended to sow discord.

## 2. Literature

The present study is amongst the first attempts by academic researchers to systematically make sense of the strategy employed by the Internet Research Agency in its "trolling" activities. In our context, extant work on Internet trolling can be divided into two categories: (i) individualistic, and (ii) coordinated. While the majority of Internet trolling research falls into the first category, a small but burgeoning stream of research is emerging with respect to the second category. An online trolling study is *individualistic* when it implicitly or explicitly assumes the Buckels et al. definition of online trolling. That is, "the practice of behaving in a deceptive, destructive or disruptive manner in a social setting on the internet with no apparent instrumental purpose" [2]. An online trolling study is coordinated if it examines strategic efforts,

undertaken by a group, to use social media platforms for subversive acts.

While the Buckels et al. definition does not explicitly rule out online trolling as a coordinated activity, studies employing this definition view online trolling as anteceded by individual nuance or as a manifestation of social practices [3]. Within the individualistic view of trolling several causal attributions have been proffered accounting for the behavior of individual trolls. In a study into the behavior of trolls on Wikipedia, [4] found that trolls repeatedly engage in intentional, harmful actions online, that they actively violate the terms of service on the platforms that they use, and that they work in isolation often with obfuscated virtual identities. That study attributes these behaviors to boredom, fun, vindictiveness and masochistic self-fulfillment. [5] identifies sexism/racism, grieving and misleading as intentional trolling behaviors on video game platforms. The individual-level antecedents to trolling reported in [5] comport with those in [4]. Another study finds a strong positive relationship between trolling behavior and a sadistic personality profile [6] while [7] find that individuals' moods predict troll behavior. Some studies focus on trolling as both individual-level and coordinated activity [3], [8].

The second theme, coordinated online trolling, is relatively nascent and heavily skewed towards Russian activities in the Occident and in Ukraine. Studies in the domain have examined retweet networks stemming from Russian IRA trolls finding that IRA trolls on twitter were vastly more popular amongst conservatives than liberals (with respect to the US) [9] that troll content did not travel far beyond Twitter, that discussion topics produced by the trolls was event-specific [10] and designed to be polarizing by tweeting in favor of "both sides" of a divisive issue [11]. In some cases, these trolls have been shown to be aggressive, targeting particular social media users and topics such that those users are dissuaded from engaging with those topics [12]. Most of the published empirical studies examining coordinated, strategic online trolling rely on data from Twitter. Russian IRU use of Facebook for trolling activities has not been examined. Further, there are no studies linking Russian IRU trolling activities with real world events. Consequently, little is known regarding the strategic deployment of trolls to: (1) exploit real world events to enhance messaging and trolling efficacy, and (2) potentially effect real world events. In addition, existing research into online trolling in general and coordinated trolling by the Russian IRA in particular has not examined the strategic timing of coordinated

troll content. To fill this gap in literature, we will analyze the relationship between Russian IRA Facebook advertisment topics and the prevailing social context around the time that those topics were broached. To do so, we construct a novel dataset that joins detailed Russian IRA Facebook advertisement data with a Black Lives Matter (BLM) protest dataset. The result is a time series consisting of both themes.

## 3. Data Collection

We noted that US House Democrats on the House Intelligence Committee released the full text as well as metadata associated with 3,393 promoted Facebook posts (advertisements). The release comprised of 3,393 individual semi-structured PDF documents. House staff appear to have scanned printed copies into PDF format; the data is embedded with the documents as images. The format of the data release, accordingly, posed a challenge for analysis. Using a Java API for Tesseract OCR (Optical Character Recognition) in concert with Oracle's PDFBox API we automated a pipeline that traverses each PDF document and applies OCR to the primary image embedded within the first page of the document.

The documents follow a standard format where the first page contains fields regarding the content timing and other metadata germane to the advertisement. The second page is the actual rendition of the advert as it would have been shown to users. We do not make use of the second page. Once an image is successfully extracted from a PDF, the Tesseract OCR API is used to extract legible text. The pipeline ends with the application of heuristics to extract data from relevant fields for structured analysis. Tesseract OCR successfully extracted 2,471 ads of the original collection of 3,393. To inspect the accuracy of the OCR process we randomly sampled 30 ads from the collection. The OCR classifier did not significantly misclassify characters in any of the sampled advertisements.

We also collected a dataset of 1,921 Black Lives Matter related protests and demonstrations from August 2014 through to May 2018. The dataset was crowd sourced on http://elephrame.com, a site dedicated to tracking these demonstrations. Below, we chart the count of BLM protests by quarter. The second half of 2016 saw a massive uptick in the count of protests.
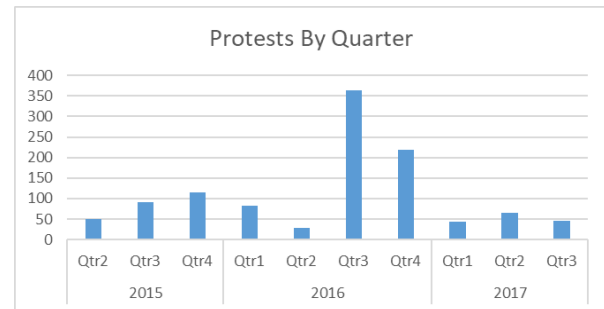

**Figure 1. Count of Protests by Quarter**

Below we chart the count of IRA advertisements by quarter. Comparing the chart below with the chart above, there is an unmistakable alignment between the protests and IRA ad buys on Facebook. Indeed, a lagged relationship. Below, the IRA ramped up its ad buying activity in the second half of 2016 in general, and in the fourth quarter in particular. Whereas BLM protests ramped up in the third quarter of 2016, the IRA appears to have acted in the following quarter, sustaining that action into 2017. Clearly, there is a need to systematically confirm this relationship.
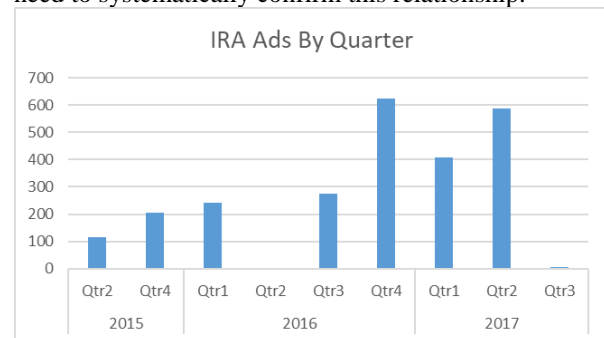

**Figure 2. Count of Ads by Quarter**

Each event recorded on http://elephrame.com is supported by a reputable news source, identified by a URL. As before, we randomly selected 30 protests from the dataset and found that all events in the sample were verifiable and correct. As the site did not provide a mechanism to download the collection, we automated the traversal of the site's pages to collect the comprehensive list. Once the comprehensive list was established, we merged both datasets into an SQL database for analysis.
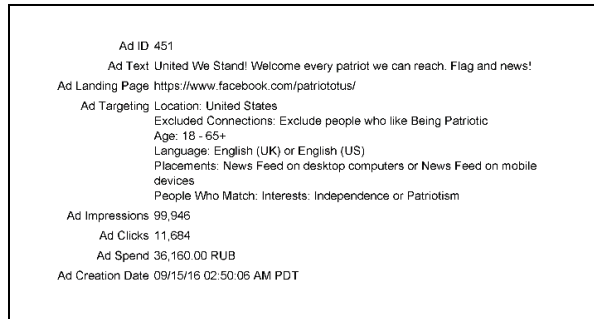
Ad ID 451
Ad Text United We Stand! Welcome every patriot we can reach. Flag and news!
Ad Landing Page https://www.facebook.com/patriototus/
Ad Targeting Location: United States
Excluded Connections: Exclude people who like Being Patriotic
Age: 18 - 65+
Language: English (UK) or English (US)
Placements: News Feed on desktop computers or News Feed on mobile devices
People Who Match: Interests: Independence or Patriotism
Ad Impressions 99,946
Ad Clicks 11,684
Ad Spend 36,160.00 RUB
Ad Creation Date 09/15/16 02:50:06 AM PDT

**Figure 3. Sample Page 1**



**Figure 4. Sample Page 2**

## 4. Topic Model

Using the non-negative matrix factorization (NMF) implementation in the Python Scikit-Learn library, we construct a topic model from the text of the advertisements. NMF is a popular decomposition technique for multivariate data where given a non-negative input matrix, $V$, two non-negative matrix factors are found, $W$ and $H$ such that their product approximates $V$ [13]. NMF begins with an $n$ x $m$ matrix, $V$, where $m$ is the number of observations in the data and $n$ is the number of "features". The $V$ matrix is factorized into an $n$ x $r$ matrix, $W$, and an $r$ x $m$ matrix, $H$, where $r$ is an input parameter specified to be lower than $m$ or $n$ [13]. In this way, NMF can be used to find low rank approximations of the input $V$ matrix.

In its application as a text mining technique, NMF can be used in document summarization or topic modelling when applied over a corpus of text documents. The textual input must be converted into a matrix representation; we elect to use a weighted representation of terms within documents by performing the term frequency dot product inverse document frequency (TF-IDF) transformation over a term frequency matrix. The result is a matrix the axes of which represent documents and terms, with fields indicating the importance of a term to a document given the distribution of that term over all documents in the corpus. The TF-IDF matrix becomes the input, $V$, to the NMF model.

After removing stop words from the text of the Russian IRA Facebook advertisements, we generate a term frequency matrix over the documents limiting the number of terms to 1000. This matrix is subsequently transformed into a TF-IDF matrix prior to non-negative matrix factorization. The number of topics to be extracted, $r$, is set to 10. The table below illustrates the results of our topic model. For each topic we list the top 10 weighted terms obtained from the $W$ factor matrix. Using these terms and analyzing the underlying advertisement text, we developed labels for the topics.

| Topic | Tokens (n=10) | Label |
|---|---|---|
| tp1 | black matters community join people proud lives blackmattersus matter life | Black Lives Matter |
| tp2 | 2nd amendment patriots lovers guns supporters defend community join priority | 2nd Amendment |
| tp3 | self defense free feel safe class event friends bring join | Self Defense Training |
| tp4 | repost women 22 powerful melanin let 30 different day choose | Misc/Unknown |
| tp5 | com https follow facebook channel twitter awww instagram www youtube | Follow Request |
| tp6 | police bm brutality man officer cop cops video officers shot | Police Brutality |
| tp7 | stop refugees news latest jobs home taking real problem islamophobia | Immigration |
| tp8 | blacklivesmatter blackandproud blackexcellence blackpride blacklove education economicempowerment | BLM Hashtags |

| | unapologeticallyafrican unapologeticallyblack knowledgeofself | |
|---|---|---|
| tp9 | veterans support brave come govspending usa learn clinton click priority | Misc/Unknown |
| tp10 | like join social issues immigration illegal agree racial looks beat | Misc/Unknown |

**Table 1. Top 10 Topic Tokens and their Labels**

## 5. Model Specification and Results

Given that existing research into online trolling in general and coordinated trolling by the Russian IRA in particular has not examined the strategic timing of coordinated troll content, this study poses the following questions: (1) are Russian IRA Facebook ads timed to coincide with polarizing events, and (2) how does the content of Russian IRA Facebook ads change to reflect strategic behavior in response to polarizing events. To glean insight into these questions, we advance the following empirical model:

$$Y_{t-i} = \beta_0 + \beta_1 Impressions_t + \sum_{j=2}^{m+2} \beta_j Topic_{tm} + \varepsilon_t$$

Where the dependent variable is a count of protests between some time $t$ and some other time $t – i$ where $i$ is denominated in days. Time periods in this model are daily intervals. In several cases, multiple IRA ads are introduced in a given day. Accordingly, impressions at time $t$ is the sum of impressions for all ads introduced during that time. A topic $m$ at time $t$ is the sum of scores for topic $m$ across all advertisements introduced at time $t$. The model coefficients are estimated using a Poisson regression to account for the non-negativity of the count based dependent variable. Impressions, the count of Facebook users who saw the advertisement is rescaled using a min-max scaler to account for the extreme dispersion in impressions; some advertisements were never seen while others were seen millions of times. In the histogram below, the sum of topic weights across the collection of advertisements is given. Topics t1 and t6, corresponding to Black Lives Matter hashtags and police brutality respectively, are clearly dominant themes in the corpus. Media reports (see above) reflect this finding, providing partial support for our topic model.

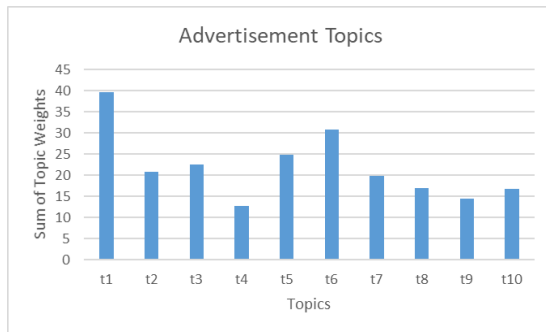| | i = | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| n =797 | **10** | | **1** | | **0** | | **-10** | | **-1** | |
| **Variables** | coef | P>|z| | coef | P>|z| | coef | P>|z| | coef | P>|z| | coef | P>|z| |
| constant | 2.601 | 0*** | 1.043 | 0*** | 0.353 | 0*** | 2.6246 | 0*** | 1.0114 | 0 |
| impressions | -0.006 | 0.962 | -0.389 | 0.224 | -0.126 | 0.764 | 0.0882 | 0.458 | -0.0274 | 0.92 |
| Black Lives Matter | -0.031 | 0.689 | -0.018 | 0.924 | 0.205 | 0.394 | 0.167 | 0.017** | 0.1614 | 0.333 |
| 2nd Amendment | 0.028 | 0.482 | 0.135 | 0.140 | 0.232 | 0.038** | 0.0372 | 0.392 | 0.1974 | 0.047** |
| Self Defense Training | -0.318 | 0.003*** | -0.965 | 0.006*** | -1.252 | 0.038** | -0.5928 | 0*** | -1.6421 | 0*** |
| Misc/Unknown | -1.189 | 0*** | -0.902 | 0.036** | -0.530 | 0.340 | -0.9672 | 0*** | -0.8205 | 0.023** |
| Follow Request | -0.547 | 0*** | -0.936 | 0*** | -1.755 | 0*** | -0.9065 | 0*** | -1.278 | 0*** |
| Police Brutality | 1.024 | 0*** | 1.084 | 0*** | 1.347 | 0*** | 0.9741 | 0*** | 1.195 | 0*** |
| Immigration | -0.016 | 0.792 | -0.133 | 0.425 | -0.007 | 0.973 | 0.1055 | 0.037** | 0.1303 | 0.234 |
| BLM Hashtags | -0.724 | 0*** | -0.290 | 0.260 | -0.998 | 0.063 | -0.6559 | 0*** | -0.7547 | 0.017** |
| Misc/Unknown | 0.235 | 0.007*** | 0.107 | 0.632 | 0.121 | 0.695 | -0.009 | 0.928 | -0.0011 | 0.996 |
| Misc/Unknown | -0.249 | 0.035** | -1.765 | 0.009*** | -3.124 | 0.005*** | -0.0625 | 0.419 | -0.09 | 0.648 |
| Log Likelihood | -6959.2 | | 2376.6 | | -1707.3 | | -7034.9 | | -2642.9 | |
| Pearson Chi-Squared | 16700 | | 3960 | | 3960 | | 16500 | | 6020 | |
| *** .01 significance level; ** .05 significance level | | | | | | | | | | |

**Table 2. Empirical Results**

**Figure 5. IRA Advertisement Topics**

Below, we present the results of our analysis. We examine different configurations of the dependent variable by toggling the size of the date range from which protest counts are computed. Recalling the model shown in the formula above, $i$ represents the outer boundaries of a date range. For instance, if $i = 30$, then, for a given time period, we are examining a dependent variable computed as the count of protests in the 30-day date range *prior* to the period. Negative values of $i$ denote a forward looking dependent variable. For instance, if $i = -30$, then, for a given time period, we are examining a dependent variable computed as the count of protests in the 30-day date range *following* the period.

When $i = 1$ or $i = -1$ the dependent variable is the count of BLM protests on and one period before a given period and, on and one period after a given period respectively. We construct our time series based on daily intervals. As such when $i = 0$ we are examining protests that occurred on the same day as the creation of a given advertisement. The results above include several significant relationships between the count of BLM protests and the topics of Russian IRA ads.

The police brutality topic is positively and significantly associated with the occurrence of BLM protests in all intervals at the .01 level. Indeed, the coefficient (to be interpreted on a logarithmic scale due to the log linking function used in a Poisson regression) values increase the closer $i$ is to 0. This may indicate that the IRA detected and responded quickly to domestic unrest regarding police violence. This finding is supported by recent work [14] showing that "Black Lives Matter protests are more likely to occur in localities where more Black people have previously been killed by police" (p. 400). It is reasonable to expect, then, that messages amplifying the killing of Black people by police will be associated with an uptick in BLM protest activity around the time of the killing. Our work does not support any causal inferences. This model cannot show that Russian IRA

advertisements related to police brutality lead to an uptick in BLM related protest activities. However, the strong association between BLM protests and police brutality related advertisements is more likely suggestive of astute opportunism on the part of the Russians. It is reasonable to assume that Russian IRA agents were monitoring Western news for reports of Black deaths stemming from police encounters.

With respect to the forward-looking models (where $i < 0$) we note positive coefficient in the police brutality topic. The strength of the positive association between these ads and BLM protests diminishes as the time interval widens. Again, this does not indicate that IRA adverts on Facebook caused protests but that upticks in protest activities as a result of police killings last for several days. We do, however, assert that this research demonstrates the need for a carefully crafted causal study into this matter.

Most topics are negatively associated with BLM protests. This finding comports with media reports about the content of the Russian IRA Facebook advertisements. If it holds that BLM protests proxy for periods of unrest, then it stands to reason that the Agency would tend to post inflammatory messages that coincide with increased protest activity. Our interpretation of these negative coefficients is as follows. During 'quiet' periods and using the count of BLM protests in a period as proxy for 'quietness,' the topic composition of IRA ads shifts away from police brutality. For instance, the self-defense training topic (tp3) is significantly negatively associated with the count of BLM protests for all tested values of $i$. This indicates that during calm periods, rather than attempt to inflame, the IRA sought to pass its Facebook accounts as normal and innocuous.

Given these results, we can proffer the following responses to our research questions. With respect to the first research question, "are Russian IRA Facebook ads timed to coincide with polarizing events" we show that Russian IRA Facebook ads are carefully timed to coincide with polarizing events. For our second research question, "how does the content of Russian IRA Facebook ads change to reflect strategic behavior in response to polarizing events", we show that the content of Russian IRA Facebook ads became increasingly inflammatory in response to polarizing events, but became mundane in response to periods of relative calm.

## 6. Limitations

There are several important limitations to this research and its findings. One such limitation is causal inference. The findings here are by no means causal in nature. This means that we cannot infer from this model that Russian IRA Facebook advertisements caused any changes in the frequency of BLM protests in subsequent periods. Another limitation is the fact that, at the time of this writing, we have been unable to extract data from about a quarter of the IRA advertisements released by the House Select Committee on Intelligence. While we do not observe any systematic bias introduced by these missing data points, we do acknowledge that the effect of the missing data is largely unknown. Finally, we acknowledge the fact that the dependent variables, where $i \neq 0$, overlap. We adjust for this by estimating model coefficients for $i = 0$ such that, for those estimates, there is no overlap. A future study will use known corrections to efficiently estimate a model under these conditions.

## 7. Conclusion

Our study shows that the Russian State attempted to sow discord in the United States via information warfare, a tactic that has gained popularity in the Russian intelligence services in recent years [1], [9], [12]. These activities are a case of coordinated online trolling. The dynamics of their use of information for subversive purposes in the United States presents an opportunity to examine general propositions about the timing of coordinated trolling activity. The use of professional trolls in one Internet Research Agency to strategically subvert normal democratic processes in the US is bourgeoning area of study. This study leverages a collection of Russian IRA produced promoted posts (ads) on Facebook released by Democrats on the US House Select Committee on Intelligence in May, 2018. Using data extracted via OCR from this release, we match advertisements to BLM protests. To do so, we take advantage of the only publicly available list of BLM protest activity. We extract a collection of 10 topics from the text of the ads. We subsequently develop a model to shed light on the association between ad topics and protest activities. We find that ads about police brutality were positively associated with protest activity by BLM activists and sympathizers. We also find that ads about all other topics were negatively associated with protest activity. These findings allow us to conclude that Russian IRA trolls were opportunistic, placing ads timed to coincide with periods of unrest and division in the United States. We can also conclude that during times of low protest activity, Russian IRA trolls turned to more innocuous topics.

This study contributes to the academic and public understanding of coordinated online troll behavior. It provides the first (to our knowledge) piece of evidence that coordinated trolling by foreign entities is a carefully timed activity. More specifically coordinated actors have a sense for the general political climate of the environments that they attack, they likely actively monitor the information landscape of that environment and alter their messaging to maximize their effect.

## 8. References

[1]     National Intelligence Council, "Assessing Russian Activities and Intentions in Recent US Elections," no. January, p. 14, 2017.
[2]     E. E. Buckels, P. D. Trapnell, and D. L. Paulhus, "Trolls just want to have fun," *Pers. Individ. Dif.*, vol. 67, pp. 97–102, 2014.
[3]     A. G. B. Cruz, Y. Seo, and M. Rex, "Trolling in online communities: A practice-based theoretical perspective," *Inf. Soc.*, vol. 34, no. 1, pp. 15–26, 2018.
[4]     P. Shachaf and N. Hara, "Beyond vandalism : Wikipedia trolls," vol. 36, no. 3, pp. 357–370, 2010.
[5]     M. D. Griffiths, "Mark D. Griffiths Adolescent trolling in online environments: A brief overview," vol. 32, no. 3, pp. 85–87, 2014.
[6]     E. E. Buckels, P. D. Trapnell, and T. Andjelovic, "Internet Trolling and Everyday Sadism: Parallel Effects on Pain Perception and Moral Judgment," *J. Pers.*, pp. 1–40, 2018.
[7]     J. Cheng, M. Bernstein, C. Danescu-niculescu-mizil, and J. Leskovec, "Anyone Can Become a Troll : Causes of Trolling Behavior in Online Discussions," 2017.
[8]     M. R. Sanfilippo, "Managing Online Trolling : From Deviant to Social and Political Trolls," pp. 1802–1811, 2017.
[9]     A. Badawy, E. Ferrara, and K. Lerman, "Analyzing the Digital Traces of Political Manipulation: The 2016 Russian Interference Twitter Campaign," 2018.
[10]     S. Zannettou, T. Caulfield, E. De Cristofaro, M. Sirivianos, G. Stringhini, and J. Blackburn, "Disinformation Warfare: Understanding State-Sponsored Trolls on Twitter and Their Influence on the Web," 2018.
[11]     L. G. Stewart, A. Arif, and K. Starbird, "Examining Trolls and Polarization with a Retweet Network," *Proc. WSDM Work. Misinformation Misbehavior Min. Web (MIS2).*, p. 6, 2018.

[12]     J. Aro, "The cyberspace war: propaganda and trolling as warfare tools," *Eur. View*, vol. 15, no. 1, pp. 121–132, 2016.

[13]     D. Lee and H. Seung, "Algorithms for non-negative matrix factorization," *Adv. Neural Inf. Process. Syst.*, no. 1, pp. 556–562, 2001.

[14]     V. Williamson, K. S. Trump, and K. L. Einstein, "Black Lives Matter: Evidence that Police-Caused Deaths Predict Protest Activity," *Perspect. Polit.*, vol. 16, no. 2, pp. 400–415, 2018.