

This article was downloaded by: [USC University of Southern California]

On: 18 October 2013, At: 13:31

Publisher: Routledge

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK



Journal of Information Technology & Politics

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/witp20>

Analyzing Political Rhetoric in Conservative and Liberal Weblogs Related to the Construction of the 'Ground Zero Mosque'

Morteza Dehghani ^a, Kenji Sagae ^b, Sonya Sachdeva ^c & Jonathan Gratch ^b

^a Brain and Creativity Institute, University of Southern California, 3620A McClintock Avenue, Los Angeles, CA, 90089, USA E-mail:

^b Institute for Creative Technologies, University of Southern California, 1205 Waterfront Dr., Playa Vista, CA, 90094-2536, USA

^c Department of Psychology, Northwestern University, 2029 Sheridan Road, Evanston, IL, 60208-2710, USA

Accepted author version posted online: 30 Jul 2013.

To cite this article: Journal of Information Technology & Politics (2013): Analyzing Political Rhetoric in Conservative and Liberal Weblogs Related to the Construction of the 'Ground Zero Mosque', Journal of Information Technology & Politics, DOI: 10.1080/19331681.2013.826613

To link to this article: <http://dx.doi.org/10.1080/19331681.2013.826613>

Disclaimer: This is a version of an unedited manuscript that has been accepted for publication. As a service to authors and researchers we are providing this version of the accepted manuscript (AM). Copyediting, typesetting, and review of the resulting proof will be undertaken on this manuscript before final publication of the Version of Record (VoR). During production and pre-press, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal relate to this version also.

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis makes every effort to ensure the accuracy of all the information (the "Content") contained in the publications on our platform. However, Taylor & Francis, our agents, and our licensors make no representations or warranties whatsoever as to the accuracy, completeness, or suitability for any purpose of the Content. Any opinions and views expressed in this publication are the opinions and views of the authors, and are not the views of or endorsed by Taylor & Francis. The accuracy of the Content should not be relied upon and should be independently verified with primary sources of information. Taylor and Francis shall not be liable for any losses, actions, claims, proceedings, demands, costs, expenses, damages, and other liabilities whatsoever or howsoever caused arising directly or indirectly in connection with, in relation to or arising out of the use of the Content.

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden. Terms & Conditions of access and use can be found at <http://www.tandfonline.com/page/terms-and-conditions>

Analyzing Political Rhetoric in Conservative and Liberal Weblogs

Related to the Construction of the ‘Ground Zero Mosque’

Morteza Dehghani¹, Kenji Sagae², Sonya Sachdeva³, Jonathan Gratch²

¹Brain and Creativity Institute, University of Southern California

3620A McClintock Avenue; Los Angeles, CA 90089, USA. mdehghan@usc.edu

²Institute for Creative Technologies, University of Southern California,

1205 Waterfront Dr., Playa Vista, CA 90094-2536, USA

³Department of Psychology, Northwestern University

2029 Sheridan Road, Evanston, IL 60208-2710, USA

Author Note

This research was supported by AFOSR Young Investigator Program and AFOSR FA9550-09-1-0507.

Morteza Dehghani is a Research Assistant Professor in Computer Science, Psychology and in the Brain and Creativity Institute at the University of Southern California. His research interests include computational social sciences and psychoinformatics.

Kenji Sagae is a Research Assistant Professor of Computer Science at University of Southern California and a Research Scientist in the USC Institute for Creative Technologies. His primary research interests are in computational linguistics and data-driven natural language processing.

Sonya Sachdeva is a postdoctoral fellow in Psychology at Northwestern University. Her research interests include cultural influences on moral cognition and environmental decision-making.

Jonathan Gratch is a Research Professor of Computer Science and Psychology at the University of Southern California (USC) and an Associate Director at USC's Institute for Creative Technologies. His research focuses on computational models of human cognitive and social processes.

ACCEPTED MANUSCRIPT

Abstract

We use different text-processing algorithms to gain insight into the political rhetoric used in conservatives and liberals weblogs. We specifically focus on the online debate regarding the issue of the ‘Ground Zero Mosque’, which has been one of the most controversial political issues in US politics in the last several years. Overall, our results show that there are significant differences in the use of various linguistic features related to sentiments of collective identity, moral concerns, and emotional dynamics between liberals and conservatives, thus, highlighting the differences between the ideological and moral frameworks of these two groups.

Keywords: blog analysis; text classification; sacred rhetoric; topic modeling; LIWC; SVM; Ground Zero mosque

There is evidence that differences in the very definitions of morality are at the root of many social-ideological differences within a country. Haidt and Graham (2007) propose that liberals and conservatives in the US have different ways of seeing the social environment around them, and rely on distinct moral structures and ideologies. Consequently, several important differences have been noted in the political rhetoric employed by these groups (Lakoff 2002, 2008; Marietta 2008, 2009). Lakoff (2008) argues that the type of language used in political discussions is of utmost importance because it “is far more than a means of expression and communication... It organizes and provides access to the system of concepts used in thinking” (p. 231). The distinct vocabulary and sets of words that each group employs result in different overarching “frames”, or metaphors, that are used to structure the complex moral issues often intertwined in politics (Lakoff, 2004). A linguistic study of presidential debates from 1976-2007 (Marietta, 2009)

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

revealed that Republicans employed sacred rhetoric, grounded in “transcendent authority and moral outrage”, more frequently and on a broader range of issues. Democrats, on the other hand, relied more on quantitative facts such as plans and projected numbers. Others have argued that the Republicans’ success in the recent elections have been largely due to their ability to find and utilize “words that work” (Luntz, 2007). In summary, the language used by liberals and conservatives in political discussions often conveys the value systems adhered to by these groups and reflects the ideological differences between them.

In this paper, we use three different statistical text-analysis techniques to examine whether such ideological differences in political views are reflected in the use of language in liberal and conservative blog posts about the Cordoba Muslim Community Center, the so-called “Ground Zero Mosque”, near the former site of the World Trade Center. We first explore whether the differences in the choice of words used by conservative and liberal bloggers is significant enough that classifiers can be trained to automatically categorize posts as conservative or liberal. We then use feature analysis to explore the most indicative features of the groups, examining what makes the posts liberal or conservative, and gaining insight into the ideologies of the groups. Using automated topic analysis (Blei, Ng & Jordan, 2003), we analyze differences in how various moral intuitions are expressed between the two groups. In line with findings of Graham, Haidt and Nosek (2009), we expect to find greater similarity in topics related to the moral domain of harm and fairness compared to other moral domains. Finally, examining further differences in language use, we use the Linguistic Inquiry and Word Count tool (Tausczik & Pennebaker, 2010) to track linguistic changes associated with affect, religiosity and sociality in the two groups over a period of one year. Our hypothesis is that if there is a greater use of sacred

ACCEPTED MANUSCRIPT

rhetoric by conservatives, it should be accompanied by an increase in the use of religious and affective words, particularly those related to anger.

Political Weblogs

In the past decade, research on political weblogs has gained significant attention. Even though only a small portion of the population actively reads weblogs, their influence on the general political atmosphere has increased significantly in the past several years. The 2004 U.S. Presidential Election was the first major political event in which parties recognized the importance of blogging as a way of connecting grassroots support, and therefore, started utilizing blogs as a major political platform (Adamic & Glance, 2005). Since then, political blog analysis has been the subject of many different lines of research¹. The aim of this paper is different from the previous lines of research, as weblogs themselves are not the subject of our research. We use the language expressed in the blogs as means of gaining better understanding of differences between the ideological and moral frameworks of conservatives and liberals. Before discussing our analysis and findings, we first outline our corpus creation method.

Corpus

On December 8, 2009, the *New York Times* published an article on plans to build the Cordoba Muslim Community center, an Islamic cultural center at a building two blocks from Ground Zero (Blumenthal & Mowjood, 2009). In response to this article, a conservative blogger criticized the project dubbing it as the “Ground Zero mosque” (Geller, 2009) and sparked a national controversy about the issue that lasted about 10 months (Elliott, 2010). The issue of the Cordoba

ACCEPTED MANUSCRIPT

Muslim Community Center, or ‘the Ground Zero Mosque’ as it came to be called, has emerged as one of the most contentious political issues in the United States in the past several years. It served to highlight the ideological differences, or the “culture war” (Hunter, 1991), between liberal and conservative moral frameworks and to a certain extent, exposed the deep prejudices that still remain toward the Muslim community. In this paper, we use the text retrieved from more than 3000 conservative and liberal blog posts related to the construction of the Ground Zero Mosque (hereafter referred to as “the Mosque”) to examine the differences in the use of linguistic features between the groups.

We had two reasons for analyzing this particular debate. First, the issue of “the Mosque” highlights the moral discrepancies in the “culture war” between liberals and conservatives. This conflict became morally significant for both groups in a short period of time and there was prevalent use of sacred rhetoric by both sides, clashing the sacred American value of religious freedom against the moral decadence of contamination of the “hallowed ground” (Davis & Dover, 2010) at Ground Zero. Second, one of the most interesting aspects of the controversy regarding “the Mosque” is the fact that it initially started on the blogosphere by a single blogger (Elliott, 2010), and most of the discussions regarding this issue took place on various different political blogs. This provided us the ability to track responses to events as they naturally unfolded, allowing longitudinal analysis of changes in different linguistic and psychological factors.

In order to compile a representative sample of the blog posts of each group, we first identified the top five most popular conservative and liberal news blogs according to the website

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

blogs.com². Next, we performed Google searches to find all posts within each of these blogs that included the word “mosque”³ and were posted between January 1, 2010 and December 31, 2010. We then automatically downloaded the HTML files (i.e. the content of the links) for all the links returned by the search queries. This included a total of 3449 blog posts, consisting of 1575 posts from the conservative blogs and 1874 from liberal blogs. Finally, we used customized scripts for each blog to remove HTML tags, headers, tables, etc. and extracted only the blog post itself and the comments on the post, ignoring all the other fields such as advertisements, blog rolls, name of the authors and dates.

Methodology

We use three different text-analysis techniques to investigate differences in the use of various linguistic features related to sentiments of collective identity, moral concerns, and emotional dynamics in liberal and conservative blog posts concerning the construction of “the mosque”. In this section we describe our methodology for analysis of these blog posts.

Group Identity

As discussed above, one of the important differences between liberals and conservatives is the type of language and rhetoric they employ in political discussions, which often reflects disparities in the ideologies and value systems of the two groups. The aim of the first analysis is to see if the differences in language that are apparent from choice of words alone are strong enough that blog posts can be automatically classified as conservative or liberal using a machine learning technique. If we are able to classify these blog posts, we then could determine the

ACCEPTED MANUSCRIPT

indicative features of each group using feature analysis and gain insight into what makes the blogs conservative or liberal.

As mentioned in the introduction, automatic political weblogs classification is an active area of research and various different techniques have been proposed for this purpose (e.g. Efron, 2004; Mullen & Malouf, 2006; Durant & Smith, 2006; Jiang & Argamon, 2008a; Jiang & Argamon, 2008b). Here, we are not interested in classifying the blogs per se, but we use classification as means of understanding the differences between the choice of words between liberals and conservatives. In a similar line of work using machine learning to examine political differences, Diermeier, Godbout, Yu and Kaufmann (2011) classified Senate speeches by first training a classifier on the speeches of the 25 most liberal and 25 most conservative senators from the 101st through 107th Congresses. They then tested their classifier on the speeches of the 25 most liberal and 25 most conservative senators of the 108th Congress and achieved an accuracy of 92%. Also, they used a similar technique to automatically classify Senate speeches by training classifiers on speeches by House representatives (Yu, Diermeier & Kaufmann, 2008). Performing a feature analysis, they reported that the most important features for Democrats included company names and words related to environmental and economic interests (e.g. Enron, ethanol, hydrogen, lakes). Conservatives, on the other hand, opted to include more words with cultural significance (e.g. cloning, unborn, abortion, marriage and homosexual).

Our approach was to use supervised machine learning, in which training data for each predefined category is needed to build a classifier. This classifier is then used to predict the category in which each new data point belongs. Support Vector Machines (SVMs), first

ACCEPTED MANUSCRIPT

introduced by Vapnik (1995), is a general learning algorithm used for binary classification. SVMs represent features, or data points, as points in space and try to find a hyperplane that is maximally distant from nearest training data points of each of the categories. In SVMs, words with the highest absolute coefficients (i.e. most positive for one group, and most negative for the other group) are considered the most informative features, and are the most indicative, or discriminative, of each category.

We used SVM^{light} (Joachims, 1999) with its default settings and the bag-of-words representation (unordered collection of words) for blog posts in this analysis. Prior to generating features vectors for classification, the documents were subjected to several pre-processing procedures. We first used a tokenizer⁴ to separate text into individual words. In order to reduce vocabulary size, we then derived⁵ word stems and mapped different forms of each word to these stems. Finally, we removed stop words, which are common words not useful for classification (e.g. “the”, “a”, “is”), along with several other categories of words such as name of the blogs and names of frequently referred to websites such as twitter.com and youtube.com. For training, we used the “term frequency–inverse document frequency” (*tf*idf*) word weighting scheme to convert documents and words in the documents to numerical document vectors. In the prediction step, given that the total number of documents is assumed to be unknown to the classifier, only word frequencies were used to represent test documents. Classification accuracy was calculated using a 10-fold cross validation, where in each run our program randomly chose a subset of the blogs from each group as the training set, and 25 other blog posts from each group as the testing sample. This process was repeated ten times and the overall accuracy of the classification was obtained by averaging over the accuracy of each of the tests.

ACCEPTED MANUSCRIPT

We also examined blog posts according to the date that they were posted. In order to get consistent number of posts per time period for both groups, we grouped blog posts into 8 consecutive time blocks (1/01-7-13, 7/14-8/10, 8/11-8/17, 8/18-8/23, 8/24-9/02, 9/03-9/13, 9/14-10/07, 10/08-12/31)⁶. In order to make sure that there would be sufficient data in each block for creating classifiers, the time blocks were chosen so that there would be at least 200 blog posts for both liberal and conservative blogs per time block. The large time blocks were necessary in order to compensate for the amount of noise existing in the files retrieved from the websites, especially the noise in the comments sections. We report the classification accuracy of our method and differences in the frequency of references to in-group vs. out-group terms in the feature vectors of the classifiers.

Moral Concerns

In the second analysis we explore similarities and differences between the two groups in how they construe different moral domains. Recent work in social psychology suggests that liberals and conservatives attend to different moral intuitions (Graham, Haidt & Nosek, 2009). Specifically, while liberals focus exclusively on the notions of harm and fairness when making moral judgments, conservatives also attend to ideas of authority, loyalty to in-group members, and purity.

We examine how the different moral intuitions of liberals and conservatives are expressed in blog posts in our corpus through the use of automatic topic modeling. We employ a hierarchical generative topic modeling technique based on Latent Dirichlet Allocation (LDA) (Blei, Ng and Jordan, 2003), which is an approach to modeling the process by which the text and the topics in a

ACCEPTED MANUSCRIPT

document are generated. LDA assumes that each document in the corpus under analysis (in our case, each individual blog post) is composed of a mixture of topics, and that each topic is a distribution over words. According to this view, documents are generated by repeatedly choosing a topic from a distribution over topics, then choosing a word from a distribution over words that corresponds to the chosen topic. Since the words in the documents are observable, the goal is to find the topic structure that generates the given collection of documents. Although LDA was originally formulated as an unsupervised topic modeling approach, where topics are induced from text without any other guidance or knowledge beyond the targeted number of topics, there are now variants that perform either unsupervised or semi-supervised modeling, where knowledge can be injected into the induced model. An interesting example of how generative unsupervised modeling can be used in the analysis of text is the Expressed Agenda model (Grimmer, 2010), where an additional level of modeling is added to reflect the author of a particular document, and therefore expose authors' priorities through the topics addressed in their language output. Similarly, Yano, Cohen and Smith (2009) used a variant of LDA to model different characteristics of a collection of political blogs gathered from November 2007 to October 2008. Their approach is to model blog posts and comments jointly, creating a framework that allows for predictions regarding comment activity associated with specific posts. They demonstrate that the main topics that emerge from their data are religion, the Iraq war, energy, domestic policy and the primary election. Hillard, Purpura & Wilkerson (2008) discuss unsupervised and supervised topic analysis techniques applied to political text, concluding that supervised modeling better serves topic classification tasks in social science research. Our approach is to use semi-supervised topic modeling, which combines the ability to allow latent

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

topics to be identified without manual supervision, while allowing expert knowledge to be taken into account to guide the topic induction process in a direction that reflects specific research issues.

To focus our analysis on moral domains, we used a variant of LDA called Latent Dirichlet Allocation with topic-in-set knowledge (Andrzejewski & Zhu, 2009) to add a level of supervision to LDA by seeding small sets of words in a subset of the topics. This semi-supervised approach combines the advantages of unsupervised topic modeling using LDA with the ability of encouraging the emergence of certain topics in the model through small sets of words selected from the outset as prior knowledge. A predefined number of topics is chosen, as in fully unsupervised LDA. However, instead of simply searching for the most probable set of latent topics, a subset of the topics can be initialized to contain specific words.

In our experiments, we set the total number of topics to 40, and used a small number of seed words (Figure 1) to initialize five topics corresponding to five foundations from Moral Foundations Theory (Graham, Haidt & Nosek, 2009). Our seed words were selected on the basis of the Moral Foundations Theory dictionary of 295 words and word stems related to each of the moral intuitions of Harm, Fairness, Authority, In-group and Purity. LDA with topic-in-set knowledge allows for the seed words to serve as soft constraints for the emerging topics; a value of zero to one is assigned to each word, which reflects the confidence with which the word belongs to its given topic. In our experiments, we assigned each of our seed words a confidence of 0.95, which is the default value used in Andrzejewski's implementation of the algorithm⁷. This level of confidence reflects high confidence that each row in Figure 1 should reflect one of the

ACCEPTED MANUSCRIPT

40 topics found by LDA, and that each of these five topics is formed around the words in each of the five rows.

In order to explore the differences in how the five moral domains are expressed between the two groups, we calculated the cosine similarity between the words for each moral intuition in each time segment between the two groups. This calculated similarity was the main dependent variable in our analysis.

<Figure 1 goes here>

Emotional Dynamics

In the last analysis, we use the Linguistic Inquiry and Word Count (LIWC) tool (Tausczik & Pennebaker, 2010) to both further investigate differences in language use between the two groups and to track linguistic changes associated with affective and social processes throughout 2010. LIWC is one of the most widely used tools for automatic text analysis in psychology, and has provided evidence for the psychological and social implications of word use in various previous studies (Pennebaker, Mehl & Niederhoffer, 2003). LIWC has also been used as a tool for tracking changes in linguistic features over time. For instance, Back, Küfner and Egloff (2010) examined the immediate negative emotional reactions on September 11, 2001 expressed in messages sent to text pagers within the US using LIWC. In a similar study, Cohn, Mehl and Pennebaker (2004) tracked psychological changes in response to the 9/11 attacks using the daily writings of 1,084 bloggers for a period of two months before and after the attacks using LIWC.

ACCEPTED MANUSCRIPT

LIWC performs word counts and catalogs words into psychologically meaningful categories (Tausczik & Pennebaker, 2010). The default LIWC2007 dictionary includes 4,500 words and word stems which define its 76 different language categories. LIWC assigns each word to specific linguistic categories, and it reports the total number of words in each category normalized by the total number of words in the document. The LIWC categories examined in this study are: social processes (e.g. talk, share, friends), affective processes (e.g. happy, cried, abandon), anger (a subcategory of affective processes) (e.g. hate, kill, annoyed) and religion (e.g. altar, church, mosque). We also created a custom Islam category which included all words in the religion category related to Islam. We use the results of LIWC, for the categories mentioned, as dependent variables in our analysis and investigate differences with regards to usage of those categories between liberals and conservatives. We also examine changes in the frequency of use words in each of the categories through different time periods.

Results

Group Identity

Overall, with a training set consisting of 750 blogs per group, our system achieved average prediction accuracy of 91.80%, $p < 0.001$. We coded the top 100 feature words with the highest absolute coefficients for each group within a classifier that achieved an accuracy of 92% for in-group and out-group membership. This coding was done relative to each subculture, for example “pelosi”, “leftist” and “socialist” were coded as out-group for conservatives, and “republican”, “right-w” and “beck” as out-group for liberals. The results show that the most important words

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

for distinguishing both conservative and liberals were those that referenced out-group members (overall: 30% out-group, 1.5% in-group, $\chi^2 = 59.083$, $p < 0.001$; Liberals: 25% out-group, 1% in-group, $\chi^2 = 23.386$, $p < 0.001$; Conservatives: 35% out-group, 2% in-group, $\chi^2 = 33.958$, $p < 0.001$) (Figure 2).

<Figure 2 goes here>

We performed the same analysis for each of the time blocks. Specifically, for each time block, our program randomly chose 175 blog posts from each group for training and another 25 posts per group for testing. Similar to the previous analysis, this process was repeated 10 times for each of the time blocks and the overall accuracy of the classifier was calculated by averaging over the 10 tests. The classification results, averaged over the 8 time periods, was 76.50% ($p < 0.001$). The accuracy of the classifier did not significantly differ between any of the time blocks. We also ran a cross-time block classification that trained classifiers on one block and tested them on the other time blocks, which obtained an accuracy of 69.43%, $p < 0.01$.

Coding the words with the highest feature weights, in classifiers which achieved accuracy closest to the mean accuracy rate of each block for in-group and out-group membership, resulted in a similar pattern as above. Within each time block, the most indicative words for both groups were references to out-group members and negative portrayals of out-group members (all p 's < 0.05).

We found that the choice of words used by bloggers belonging to these two ideological groups were, indeed, distinct enough that our system could classify their posts as conservative or liberal with an accuracy of 91.80%. Even though we expected that this difference would diminish for

ACCEPTED MANUSCRIPT

posts within each time block, due to similar topics of discussion and reduced training data per time block, we were able to classify blog posts within each block with an average accuracy of 76.50%. The last classification result for across time-block classification was 69.43%. Although, this accuracy is lower than previous ones, it indicates an important finding about the choice of words used by the two groups. This result indicates that we can reliably distinguish between the rhetoric used by liberals and conservatives, even if we construct classifiers from the text of one time block and test it on different time blocks. This ability to perform cross-time block classification supports the notion that these two groups differ not only in how they discuss an event, but also in the type of rhetoric used across various topics of discussion. Feature analysis revealed that the most distinctive aspect of either liberal or conservative blogs was not the description, or the ideology, of the in-group, but rather the use of words related to the negative portrayal of the out-group.

Moral Concerns

A 5 (morality domain: Harm, Fairness, Authority, Ingroup, Purity)* 2(Virtue/Vice) ANOVA revealed an overall main effect of moral domain $F(4, 70) = 4.913, p = 0.001$, and main effect of Virtue/Vice $F(1, 70), p=0.005$ (Figure 4). Overall, there was significantly more similarity in Virtue domains than in Vice domains $t(78) = 2.623, p = 0.010$. There was also a significant difference between similarity in Harm domains compared to Fairness domains $t(30) = 2.692, p = 0.011$, with higher similarity in Harm. Also, there was higher similarity in Harm domains compared to Authority $t(30) = 2.1381, p = 0.04$ and Purity $t(30) = 2.389, p = 0.02$. Finally, there

was higher similarity in In-group value compared to Fairness $t(30) = 3.246, p = 0.003$ and Purity $t(30) = 2.427, p = 0.02$.

<Figure 4 goes here>

Consistent with the findings of Haidt and colleagues (Graham, Haidt & Nosek 2009; Haidt & Graham 2007), our analysis demonstrates that there are significant differences in how liberals and conservatives construct their moral belief systems. Specifically, our analysis shows that the greatest degree of similarity between the moral intuitions of liberals and conservatives is within the domain of Harm (Figure 4). This finding is consistent with Haidt et al.'s (2009; 2007) assertion that while Harm is universally accessible, liberals and conservatives differ on the extent to which they rely on Purity, Authority and Ingroup. However, contrary to Haidt et al., we find very low similarity in the topic words related to the domain of Fairness. Further work could examine whether this difference is due to what constitutes Fairness for liberals and conservatives. One piece of evidence that might support this perspective is the observed main effect of Virtue versus Vice terms. We observed greater similarity between conservative and liberal rhetoric for Virtue terms than Vice terms in the domain of Harm, Fairness and Authority but not for Purity and Loyalty. These results indicate that the positive aspect of fairness show less cross-cultural variability than the negative aspects (i.e. fairness, equity vs. inequality, injustice). We suggest that injustice might be a different construct for liberals than for conservatives. Another reason for the difference in the Fairness domain is that significantly more numbers and digits make up the Vice aspect of Fairness for liberals than for conservatives (To demonstrate this visually, in Figure 5, we replaced all digits and numbers in Vice Fairness, for

both liberals and conservatives, with the word “number”). This is consistent with the findings of Marietta (2009) who also noted that liberals tend to rely more on quantitative facts and projected numbers. Albeit, we show that this phenomenon is only relevant within the Vice aspect of Fairness.

<Figure 5 goes here>

Emotional Dynamics

First, we examined how linguistic features for affective processes changed over time. A repeated measures ANOVA, where the first factor was time and the second factor was group, determined an overall main effect of time for affect $F(7, 56) = 3.600, p = 0.003$. The same test revealed that the interaction between time and groups approached significance $F(7, 56) = 2.331, p = 0.037$. There was a positive correlation between time and use affective words for conservatives in the first six time blocks, leading up to and including the 9/11 time block $r = 0.852, p = 0.031$, however this correlation did not reach significant for liberals $r = 0.194, p = 0.712$. For conservatives, there was a significant increase in the use of affective words between the first time block and the 9/03-9/13 time block $t(8) = 4.462, p = 0.002$. Furthermore, in the 9/05-9/13 block, the differences in this category between the two groups became significant $t(8)=2.4256, p = 0.0415$. Overall, the amount of affective words used by conservative websites was higher than liberal websites $t(78) = 2.360, p = 0.021$ (Figure 6).

A repeated measures ANOVA with a Greenhouse-Geisser correction, where the first factor was time and the second factor group, determined an overall main effect of time for anger $F(2.503,$

ACCEPTED MANUSCRIPT

20.024) = 4.773, $p = 0.015$. The same test revealed that a significant interaction between time and groups $F(2.503, 20.024) = 3.458, p = 0.042$. Even though the overall difference in the use of words related to anger between the two groups did not reach significance $t(78) = 1.104, p = 0.273$, this difference became significant at the 9/05-9/13 time block $t(8) = 3.9544, p = 0.004$.

As shown in the graphs, there was a sharp decrease in the use of affective words and anger in the last time block, especially for conservatives, which is an indication of these processes returning to baseline rates (there was no significant difference between the first and last time blocks in any of the emotion categories mentioned above for either of the two groups).

Another repeated measures ANOVA was ran for the religion category. There was a main effect of time $F(7, 56) = 5.462, p < 0.001$, and the interaction between time and group approached significance $F(7, 56) = 1.914, p = 0.084$. For conservatives, there was a positive correlation between anger and religion $r = 0.7204, p = 0.044$. Correlating the Islam sub-category with anger indicated that the correlation between anger and religion was not due to use of words related to Islam ($r = 0.6406, p = 0.0870$), as the correlation of anger was stronger with the sub-category of religion not including words related to Islam $r = 0.770, p = 0.025$.

A repeated measures ANOVA with a Greenhouse-Geisser correction for social processes revealed an overall main effect of time $F(7, 56) = 2.539, p = 0.024$. For both groups, there was increase in social orientation over time (conservatives: $r = 0.909, p = 0.05$; liberals: $r = 0.780, p = 0.040$) which, similar to other factors, returned to baseline after the 9/11 timeblock. Also the use of words related to social processes was higher for conservatives than for liberals $t(78) = 4.1151, p < 0.001$.

ACCEPTED MANUSCRIPT

<Figure 6 goes here>

Overall, there were significant differences in the use of words related to affective and social processing between conservatives and liberals. As our results show, for conservatives there was a significant increase in the use of words related to affect, and anger, in periods leading up to the anniversary of 9/11. These changes in the choice of words used in the posts reflect underlying differences in the type of rhetoric employed, and subsequent changes in emotional responses.

For conservatives the rise in the use of words related to anger was positively correlated with the use of religious words, which is an indication of an increased reliance on sacred rhetoric. The use of sacred rhetoric has been linked to the emergence of sacred values (Marietta, 2008; Dehghani, Atran, Iliev, Sachdeva, Medin & Ginges, 2009; Dehghani, Iliev, Sachdeva, Atran, Ginges & Medin, 2010), as values that get tied to religion more easily achieve a sacred status (Marietta, 2009). Previous work has shown violations of sacred values result in anger and moral outrage (e.g. Tetlock, 2003; Ginges, Atran, Medin, & Shikaki, 2007).

Traumatic and upsetting events are often followed by an increase in social processes such as seeking of social support, increase in collective orientation and social sharing (Mehl & Pennebaker, 2003). Our results indicate that there were increases in the use of words related to social processing by both conservatives and liberals over time, which may have been due to attempts to validate their threatened cultural worldview (Pyszczynski, Greenberg, Solomon, Arndt & Schimel, 2004), and to facilitate social sharing (Rimé, Finkenauer, Luminet, Zech, & Philippot, 1998).

General Discussion and Conclusions

In this paper we analyzed conservative and liberal blogs posts related to the “Ground Zero” mosque in order to examine and gain better understanding of differences between the ideological and moral frameworks of conservatives and liberals reflected in linguistic characteristics. Using three different statistical text analysis techniques, we demonstrated that there are significant differences between liberals and conservatives in the use of various linguistic features and the choice of words.

In the first analysis, we used a machine learning technique to both automatically classify the blogs based on the group they were written by, and to examine the indicative features which made these blogs liberal or conservative. Our results indicate that words that reference out-group members and are used for out-group derogation are most characteristic of the ideology of a group (whether liberal or conservative). This indicates that, at least in political debates, the values and ideas that contribute to the identity of each group predominantly consist of stereotypes of the out-group.

To further examine the differences between the moral value systems adhered to by each group, we used a non-parametric Bayesian topic modeling approach to enable the unsupervised detection of topics in our corpus of liberal and conservative blogs. We used a small set of words selected from the Moral Foundations dictionary (Graham, Haidt & Nosek 2009) as seeds to encourage the emergence of topics related to different moral domains. Interestingly, although many of the same topics were discussed in liberal and conservative blog posts, our study yielded significant differences in how moral intuition terms were used by the two groups.

ACCEPTED MANUSCRIPT

By examining posts in different time blocks, we also showed that there was an increase in words related to affective processes and anger over time, especially for conservatives. This increase is likely related to the use of sacred rhetoric, as there was a significant correlation between anger and the use of religious words. In the third analysis, by examining posts in different time blocks, we were able to demonstrate the increase in words related to affective processes and anger over time, especially for conservatives. This increase may also be related to the use of sacred rhetoric, as reflected in the significant correlation between anger and the use of religious words.

In conclusion, by analyzing more than 3000 conservative and liberal blog posts related to the constructions of the “Ground-Zero Mosque”, our results confirm significant differences between liberals and conservatives in the use of language related to sentiments of collective identity, moral concerns and the emotions they elicit. We believe that the ability to perform this type of mass text analysis and to track changes in different psychological processes over different periods of time, as they naturally unfold among diverse cultural groups, can provide new insights which arguably cannot be achieved in an experimental setting inside the lab.

References

Ackland, R. (2005). Mapping the U.S. political blogosphere: Are conservative bloggers more prominent? Australian National University.

Adamic, L. A., & Glance, N. (2005). The political blogosphere and the 2004 US election: Divided they blog. *In Proceedings of the 3rd International Workshop on Link Discovery*, 36-43. Chicago, IL

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Andrzejewski, D. & Zhu, X. (2009) Latent Dirichlet Allocation with Topic-in-Set Knowledge. *Proceedings of the NAACL 2009 Workshop on Semi-supervised Learning for NLP*. Association for Computational Linguistics, Stroudsburg, PA, USA, 43-48.

Back, M. D., Küfner, A. C. P., & Egloff, B. (2010). The emotional timeline of September 11, 2001. *Psychological Science* 21: 1417-1419.

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research (JMLR)*, 3, Mar. 2003, 993-1022.

Blumenthal, R. & Mowjood, S. (December 8, 2009). Muslim prayers and renewal near Ground Zero. *New York Times*. Available at: <http://www.nytimes.com/2009/12/09/nyregion/09mosque.html>

Cohn, M. A., Mehl, M. R., & Pennebaker, J. W. (2004). Linguistic markers of psychological change after September 11, 2001. *Psychological Science*, 15, 687-693

Davis, L. & Dover, E. (Aug 22, 2010). Ground Zero Mosque opponents, supporters turn out to demonstrate. *ABC News*. Available at: <http://abcnews.go.com/US/Politics/ground-mosque-opponents-supporters-turn-demonstrate/story?id=11455698>

Dehghani, M., Iliev, R., Sachdeva, S., Atran, S., Ginges, J. & Medin, D. (2009). Emerging sacred values: Iran's nuclear program. *Judgment and Decision Making*, 4, 7, 930-933.

Dehghani, M., Atran, S., Iliev, R., Sachdeva, S., Medin, D. & Ginges, J. (2010). Sacred values and conflict over Iran's nuclear program. *Judgment and Decision Making*, 5, 7, 540-546.

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Diermeier, D., Godbout, J. F., Yu B., & Kaufmann, S. (2011). Language and ideology in Congress. *British Journal of Political Science*.

Durant, K. T. & Smith, M. D. (2006). Mining sentiment classification from political web logs. *In Proceedings of Web Mining and Web Usage Analysis '06*.

Efron, M. (2004). Cultural orientation: Classifying subjective documents by cociation analysis. *AAAI Fall Symposium on Style and Meaning in Language, Art, and Music*.

Elliott, J. (Aug 16, 2010). How the "ground zero mosque" fear mongering began. *Salon.com*. Available at: http://www.salon.com/news/politics/war_room/2010/08/16/ground_zero_mosque_origins

Farrell, H., & Drezner D.W. (2008). The power and politics of blogs. *Public Choice*. 134(1-2): 15-30.

Geller, P. (Dec 8, 2009). Giving thanks. Available at: http://atlashrugs2000.typepad.com/atlas_shrugs/2009/12/giving-thanks.html

Gentzkow, M. & Shapiro, J. (2010). Ideological Segregation Online and Offline, *National Bureau of Economics Research Working Paper* No. 15916.

Ginges, J., Atran, S., Medin, D. & Shikaki, K. (2007). Sacred bounds on rational resolution of violent political conflict. *Proceedings of the National Academy of Sciences*, 104, 7357-7360.

Graham, J., Haidt, J., & Nosek, B. A. (2009). Liberals and conservatives rely on different sets of moral foundations. *Journal of Personality and Social Psychology*, 96, 1029-1046

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Grimmer, J. (2010). A Bayesian Hierarchical Topic Model for Political Texts: Measuring Expressed Agendas in Senate Press Releases. *Political Analysis*, 18(1), 1-35.

Gumbrecht, M. (2004). Blogs as “protected space”. In *WWW 2004 Workshop on the Weblogging Ecosystem: Aggregation, Analysis and Dynamics, 2004, at WWW'04: the 13th international conference on World Wide Web*.

Haidt, J., & Graham, J. (2007). When morality opposes justice: Conservatives have moral intuitions that Liberals may not recognize. *Social Justice Research*, 20, 98-116

Hargittai, E., Gallo, J., & Kane, M. (2007). Cross ideological discussions among conservative and liberal bloggers. *Public Choice* 134(1): 67–86.

Hillard, D., Purpura, S., & Wilkerson, J. (2008). Computer assisted topic classification for mixed methods social science research. *Journal of Information Technology and Politics*, 4(4): 31–46.

Hunter, J. D. (1991). *Cultural wars: The struggle to define America*. New York: Basic Books.

Jiang, M., & Argamon, S. (2008). Finding Political Blogs and Their Political Leanings. Paper presented at the *Text Mining 2008, Workshop at the SIAM International Conference on Data Mining*.

Jiang, M., & Argamon, S.. (2008). Political Leaning Categorization by Exploring Subjectivities in Political Blogs. In *the Proceedings of the Fourth International Conference on Data Mining, DMIN'08*.

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Joachims, T. (1999). Making Large-Scale SVM learning practical. In: *Advances in Kernel methods - Support Vector Learning*, B. Schölkopf, C. Burges, and A. Smola (ed.). MIT Press, Cambridge, MA.

Karpf, D. (2008a). Measuring influence in the political blogosphere. In *Politics and Technology Review, George Washington University's Institute for Politics, Democracy & the Internet*. 33–41.

Karpf, D. (2008b). Understanding blogspace. *Journal of Information Technology and Politics*, 5(4), 369–385.

Krishnamurthy, S. (2002). The multidimensionality of blog conversations: The virtual enactment of September 11. In *Internet Research 3.0*

Lakoff, G. (2002). *Moral politics: How Liberals and Conservatives think*. University of Chicago Press, Chicago, IL.

Lakeoff, G. (2004). *Don't think of an elephant! Know your values and frame the debate. The essential guide for progressives*. White River Junction, VT: Chelsea Green.

Lakoff, G. (2008). *The political mind: Why you can't understand 21st-century politics with an 18th-century brain*. Viking, New York, NY.

Landauer, T. K. & Dumais, S. T. (1997), A Solution to Plato's Problem: The Latent Semantic Analysis Theory of Aquisition, Induction and Representation of Knowledge, *Psychological review*, 104 (2), 211-240.

Luntz, F. (2007). *Words that work*. New York: Hyperion.

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Marietta, M. (2008) "From my cold, dead hands": Democratic consequences of sacred rhetoric. *Journal of Politics*. 70, 3:767-779.

Marietta, M. (2009). The absolutist advantage: sacred rhetoric in contemporary presidential debate. *Political Communication*. 26, 4:388-411

Mehl, M. R. & Pennebaker, J. W. (2003). The social dynamics of a cultural upheaval: Social interactions surrounding September 11, 2001. *Psychological Science*, 14, 579-585.

Mishne, G., & Glance, N. (2006, May). Leave a reply: An analysis of weblog comments. Paper presented at the workshop *Weblogging Ecosystem: Aggregation, Analysis and Dynamics*, Edinburgh, UK.

Mullen, T. & Malouf, R. (2006). A preliminary investigation into sentiment analysis of informal political discourse. In *Proceedings of the AAAI-2006 Spring Symposium on Computational Approaches to Analyzing Weblogs*, 159 – 162.

Pennebaker, J. W., Mehl, M.R., & Niederhoffer, K. (2003). Psychological aspects of natural language use: Our words, our selves. *Annual Review of Psychology*, 54, 547-577.

Pyszczynski, T., Greenberg, J., Solomon, S., Arndt, J., & Schimel, J. (2004). Why do people need self-esteem? A theoretical and empirical review. *Psychological Bulletin*, 130, 435-468.

Rimé, B., Finkenauer, C., Luminet, O., Zech, E., & Philippot, P. (1998). Social sharing of emotion: New evidence and new questions. In W. Stroebe & M. Hewstone (Eds.), *European Review of Social Psychology* (Vol. 9, pp.145-189). Wile, Chichester.

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Tausczik, Y., & Pennebaker, J. W. (2010). The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology, 29*, 24-54.

Trammell, K. D. (2006). Blog offensive: An exploratory analysis of attacks published on campaign blog posts from a political public relations perspective. *Public Relations Review, 32*(4), 402-406.

Tetlock, P. (2003). Thinking the unthinkable: sacred values and taboo cognitions. *Trends in Cognitive Sciences, 7*, 320–24.

Trevino, E. M. (2005). Blogger motivations: Power, pull, and positive feedback. *In Internet Research 6.0*.

Vapnik, V. N. (1995) *The Nature of Statistical Learning Theory*. Springer, New York.

Wallsten, K. J. (2006). Agenda Setting and the Blogosphere: An Analysis of the Relationship between Mainstream Media and Political Blogs. *Review of Policy Research, 24*(6): 567-587.

Wallsten, K. (2007). Political blogs: Transmission belts, soapboxes, mobilizers, or conversation starters? *Journal of Information Technology and Politics, 4*(3), 19–40.

Wallsten, K. (2010). “Yes We Can.”: How online viewership, blog discussion, campaign statements, and mainstream media coverage produced a viral video phenomenon. *Journal of Information Technology & Politics, 7*(2–3), 163–181.

Williams, A. P., Trammell, K. D., Postelnicu, M., Landerville, K. D., & Martin, J. D. (2005). Blogging and hyperlinking: Use of the web to enhance viability during the 2004 US campaign. *Journalism Studies, 6*(2), 177–186.

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Yano, T., Cohen, W. W., & Smith, N. A. (2009). Predicting response to political blog posts with topic models. *In Proceedings of NAACL*, 477–485.

Yu, B., Diermeier, D. & Kaufmann, S. (2008). Classifying party affiliation from political speech. *Journal of Information Technology in Politics*, 5, 33-48.

Footnotes

¹ These include: examination of cross-linkage by political affiliation (Farrell & Drezner, 2008; Adamic & Glance, 2005; Ackland, 2005; Hargittai, Gallo & Kane, 2008; Gentzkow & Shapiro, 2010), analysis of different blogging activities and their effects on readers (Karpf 2008a; Karpf, 2008b; Wallsten, 2007), role of blogs in election campaigns (e.g. Williams, Trammell, Postelnicu, Landerville & Martin, 2005), impact of political blogs on mainstream news coverage (e.g. Wallsten, 2007 & 2010), analysis of comment section of blogs (e.g. Mishne & Glance, 2006; Trevino, 2005; Gumbrecht, 2004; Krishnamurthy, 2002; Trammell, 2006), and automatic party classification (determining whether a blogpost is from a conservative blog or liberal blog) (e.g. Efron, 2004; Mullen & Malouf, 2006; Durant & Smith, 2006; Jiang & Argamon, 2008a; Jiang & Argamon, 2008b).

²The conservative blogs we chose for our analysis are the following: hotair.com, reason.com, redstate.com, rightwingnews.com and townhall.com, and the liberal blogs are: crooksandliars.com, dailykos.com, huffingtonpost.com, thinkprogress.com and wonkette.com.

³ We tried several search words, and the “mosque” appeared to be the best keyword for finding posts related to the Ground-Zero Mosque debate

ACCEPTED MANUSCRIPT

⁴ For tokenization, we used the Word Splitter tool, available at http://cogcomp.cs.illinois.edu/page/tools_view/8

⁵ To derive word stems, we used the lisp implementation of the Porter stemmer, available at <http://tartarus.org/~martin/PorterStemmer/>

⁶ There was no significant difference in word count of blog posts between the liberal and conservative groups and in any of the time periods.

⁷ http://pages.cs.wisc.edu/~andrzej/research/zl_lda.html

ACCEPTED MANUSCRIPT

Figure 1: Seed words used to generate topics related to different moral intuitions

Domain	Virtue	Vice
Harm	safe defend protect	harm war kill
Fairness	fair equal justice	unfair unequal unjust
Ingroup	together nation family	foreign enemy terrorist
Authority	duty law honor	rebel betray traitor
Purity	sacred preserve pure	disgust sin disease

ACCEPTED MANUSCRIPT

Figure 2: Top 20 words with the highest weight for each group within a classifier that achieved 92% accuracy. Words are listed in decreasing weight order. All words were converted into lower case, and in order to reduce vocabulary size, word stems were used in classification.

Conservative: liber, leftist, govern, koran, august, obamacar, obama, left, lefti, conserve, law, murder, rino, infidel, union, allah, gzm, lib, properti, stan

Liberal: center, gingrich, commun, republican, fox, religi, corpor, beck, tenesse, muslin, park, krauthamm, adl, bigotri, fear, filibust, right-w, wingnut, terminu

ACCEPTED MANUSCRIPT

ACCEPTED MANUSCRIPT

Figure 3: Topic Clusters emerged using the seed words

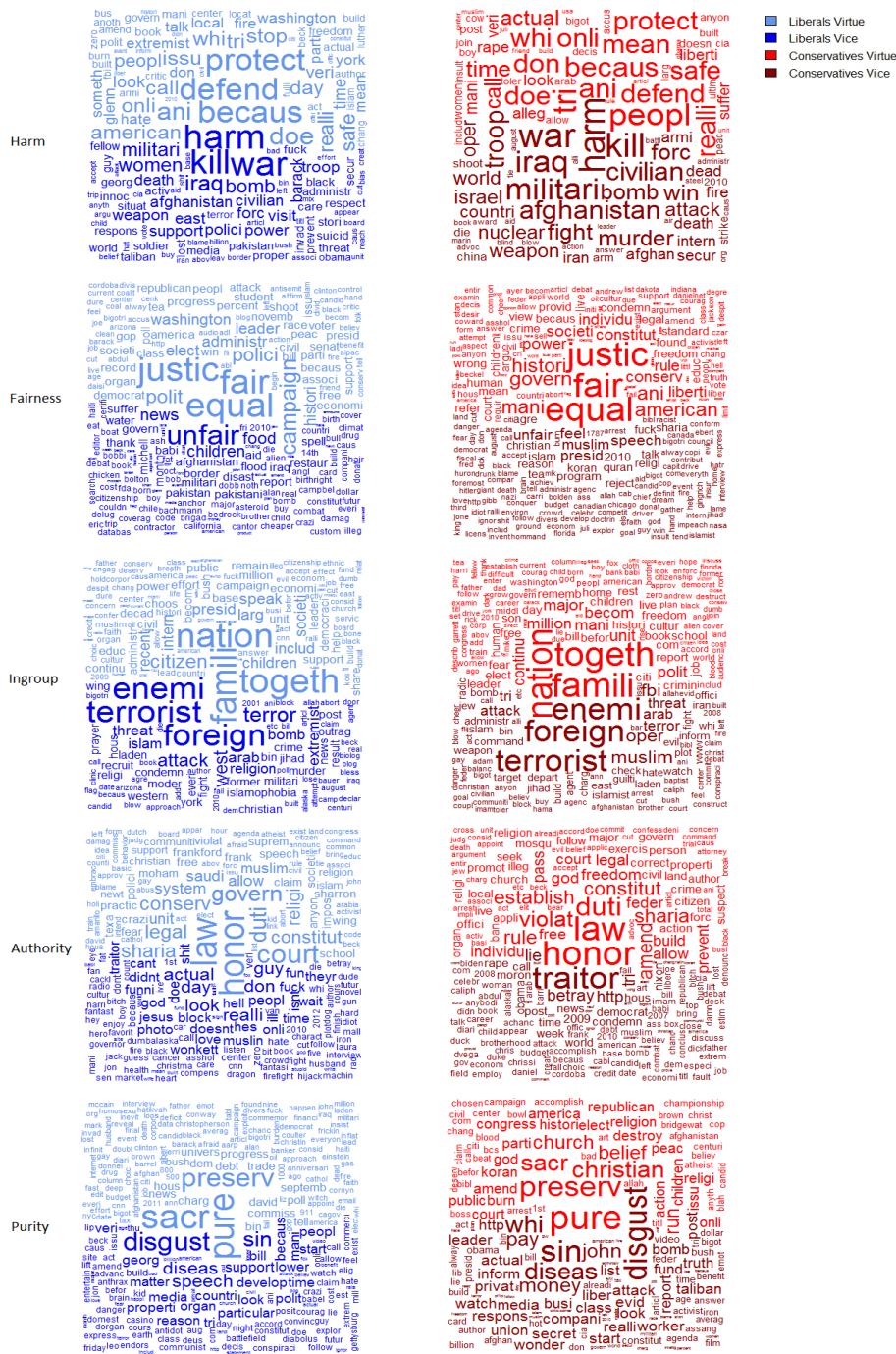
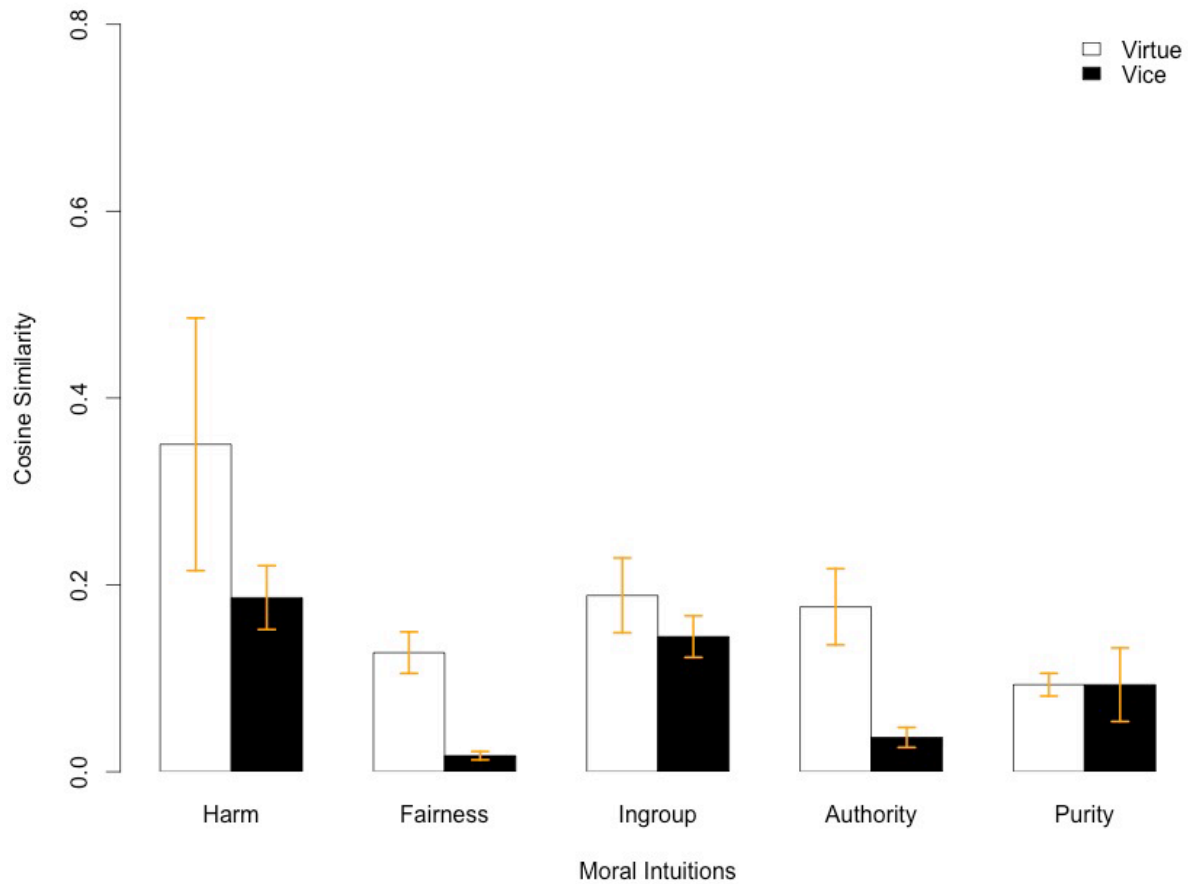


Figure 4: Cosine similarity between the moral topics emerged from conservative and liberal blogs



ACCEPTED MANUSCRIPT

Figure 5: To demonstrate the prevalent use of numbers and figures in the vice aspect of the domain of Fairness by liberals, all digits and numbers were replaced with the word “number” in Fairness vice and virtue for both groups



ACCEPTED MANUSCRIPT

Figure 6: Changes in different psychological processes captured by LIWC

