

# Combining natural language processing and network analysis to examine how advocacy organizations stimulate conversation on social media

Christopher Andrew Bail<sup>a,b,1</sup>

<sup>a</sup>Department of Sociology, Duke University, Durham, NC 27708; and <sup>b</sup>Network Analysis Center, Duke University, Durham, NC 27708

Edited by Peter S. Bearman, Columbia University, New York, NY, and approved August 5, 2016 (received for review May 6, 2016)

Social media sites are rapidly becoming one of the most important forums for public deliberation about advocacy issues. However, social scientists have not explained why some advocacy organizations produce social media messages that inspire far-ranging conversation among social media users, whereas the vast majority of them receive little or no attention. I argue that advocacy organizations are more likely to inspire comments from new social media audiences if they create “cultural bridges,” or produce messages that combine conversational themes within an advocacy field that are seldom discussed together. I use natural language processing, network analysis, and a social media application to analyze how cultural bridges shaped public discourse about autism spectrum disorders on Facebook over the course of 1.5 years, controlling for various characteristics of advocacy organizations, their social media audiences, and the broader social context in which they interact. I show that organizations that create substantial cultural bridges provoke 2.52 times more comments about their messages from new social media users than those that do not, controlling for these factors. This study thus offers a theory of cultural messaging and public deliberation and computational techniques for text analysis and application-based survey research.

computational social science | advocacy organizations | public deliberation | networks | natural language processing

During the last half century, autism spectrum disorder (ASD) diagnoses have increased dramatically within the United States (1). In the absence of scientific consensus about the root causes of ASDs, a pressing public conversation has begun among a large group of advocacy organizations about the root causes of ASDs. These organizations promote a range of different explanations for the rise in ASDs, from epigenetic factors to a heavily discredited narrative that links these disorders to routinely administered childhood vaccines (2, 3). The trajectory of this public conversation has critical implications for the estimated \$236 billion spent annually on research, treatment, and support for the 1 in 68 children diagnosed with this disorder each year (4).

Conversation about contentious issues such as the cause of ASDs is often described as the soul of democracy (5). This is not only because public conversations spread awareness about advocacy issues but also because they often provoke deliberation about how they should be addressed. Until recently, advocacy organizations were forced to circumnavigate media gatekeepers to stimulate public conversation about their cause. However, the advent of social media offers all advocacy organizations the potential to stimulate far-ranging conversations that spread rapidly across diverse groups of people (6).

Despite the promise of social media to democratize public deliberation about advocacy issues, Facebook and other social media sites have also increased the sheer scale of messages that compete for public attention each day (7, 8). Although some advocacy organizations continue to stimulate large conversations about their cause, the vast majority receive little or no attention. Social scientists have not yet produced a theory of how advocacy organizations stimulate public conversation from new audiences on social media

sites, despite the rapidly increasing influence of these forums on the trajectory of public debate about advocacy issues.

## Cultural Networks and Bridges

One of many important factors that may determine whether advocacy organizations stimulate large social media conversations is how the content of their messages fits into preexisting discourse about an advocacy issue. Although social network analysis is typically used to describe friendships or other relationships between individuals, it can also be used to describe relationships between actors via the types of messages or ideas they produce (9–14). Fig. 1 is an example of a small region within such a “cultural network.” Each node describes an actor engaged in public conversation about an advocacy issue, and the edges between the nodes represent those who are discussing similar issues within the social media advocacy field. One of the clusters of social media users pictured at  $t_1$  in Fig. 1 might be discussing the relationship between vaccines and ASDs, for example, whereas the other might be composed of those assessing epigenetic factors.

Advocacy organizations that aim to stimulate public conversation about ASDs might craft messages that address social media users within either of these clusters. However, messages that target individual discursive clusters within a social media advocacy field will have finite appeal that is proportional to the size of each cluster. What is more, organizations that produce messages about themes that are already well-established risk appearing redundant or otherwise unremarkable. In contrast, organizations that produce messages about entirely new discursive themes may be similarly ignored because they are perceived as irrelevant to the

## Significance

Social media provide potent opportunities for advocacy organizations to shape public debates because of the rapidly increasing number of people who frequent such forums each day. However, social scientists have not yet explained why some advocacy organizations create large-scale public debate whereas others do not. Using automated text analysis and a Facebook application, I found that advocacy organizations are more likely to stimulate conversation if they produce messages that link discursive themes that are usually discussed in isolation from each other. Such messages not only resonate with multiple audiences, but also put such audiences in conversation with each other. This manuscript thereby contributes a theory of public deliberation on social media for the emerging field of computational social science.

Author contributions: C.A.B. designed research, performed research, contributed new reagents/analytic tools, analyzed data, and wrote the paper.

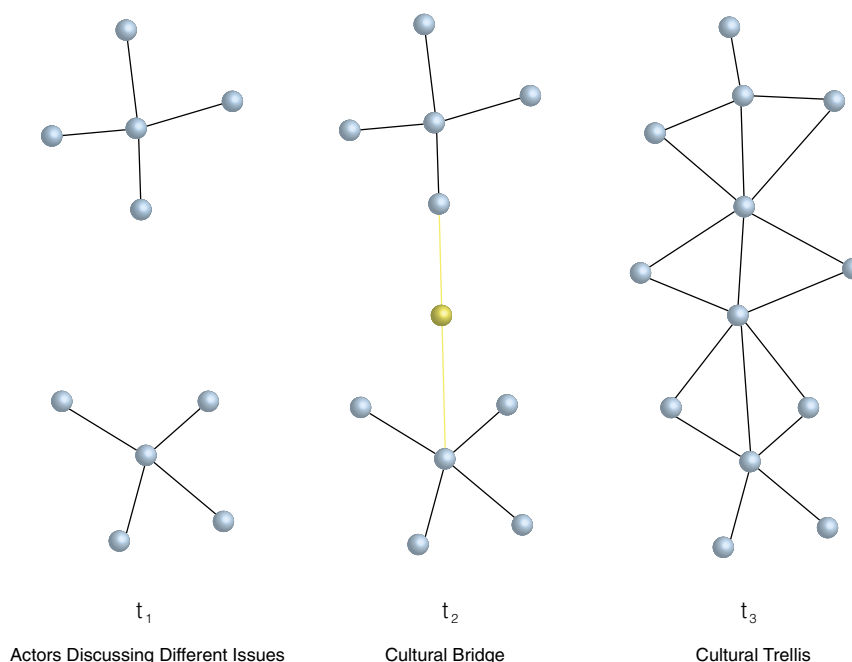
The author declares no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

<sup>1</sup>Email: christopher.bail@duke.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1607151113/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1607151113/-DCSupplemental).



**Fig. 1.** Hypothetical cultural network in which nodes represent actors engaged in conversation about an advocacy issue and edges between them describe similarities in the content of their messages. I argue that advocacy organizations are most likely to stimulate comments from new social media audiences if they create “cultural bridges,” or produce messages that connect discursive themes that are seldom discussed together. Such messages may not only provoke comments from multiple audiences but also put these audiences into conversations with one another, creating new, hybrid conversational themes, or “cultural trellises,” within a social media advocacy field.

central concerns of those within an advocacy field. The ideal message, in other words, must be both new and familiar.

I argue that advocacy organizations can achieve an optimal blend of novelty and familiarity by creating “cultural bridges,” or messages that link discursive themes within an advocacy field that are seldom discussed together. For example, the organization represented by the yellow node at  $t_2$  in Fig. 1 could create a cultural bridge by producing a social media message about the lack of scientific evidence that links ASDs to either vaccines or epigenetic factors. Such messages may stimulate comments from social media users who regularly discuss vaccines or epigenetics because they contain discursive themes that may appear new and familiar from the perspective of both populations within the social media advocacy field.

Advocacy organizations that create cultural bridges not only expand the universe of social media users who might engage with their message but also puts them into conversation with each other. The aforementioned message about the lack of scientific evidence that links ASDs to vaccines or epigenetic factors, for example, might inspire argument between those within each cluster. Or, interaction between those in multiple clusters might create what I call a “cultural trellis,” or a hybrid conversational theme such as the concept of neurodiversity (the belief that ASDs result from natural neurological variation). As  $t_3$  in Fig. 1 shows, such new themes can emerge because of repeated interactions between those in discursive clusters that are typically isolated from each other.

### Analytic Approach

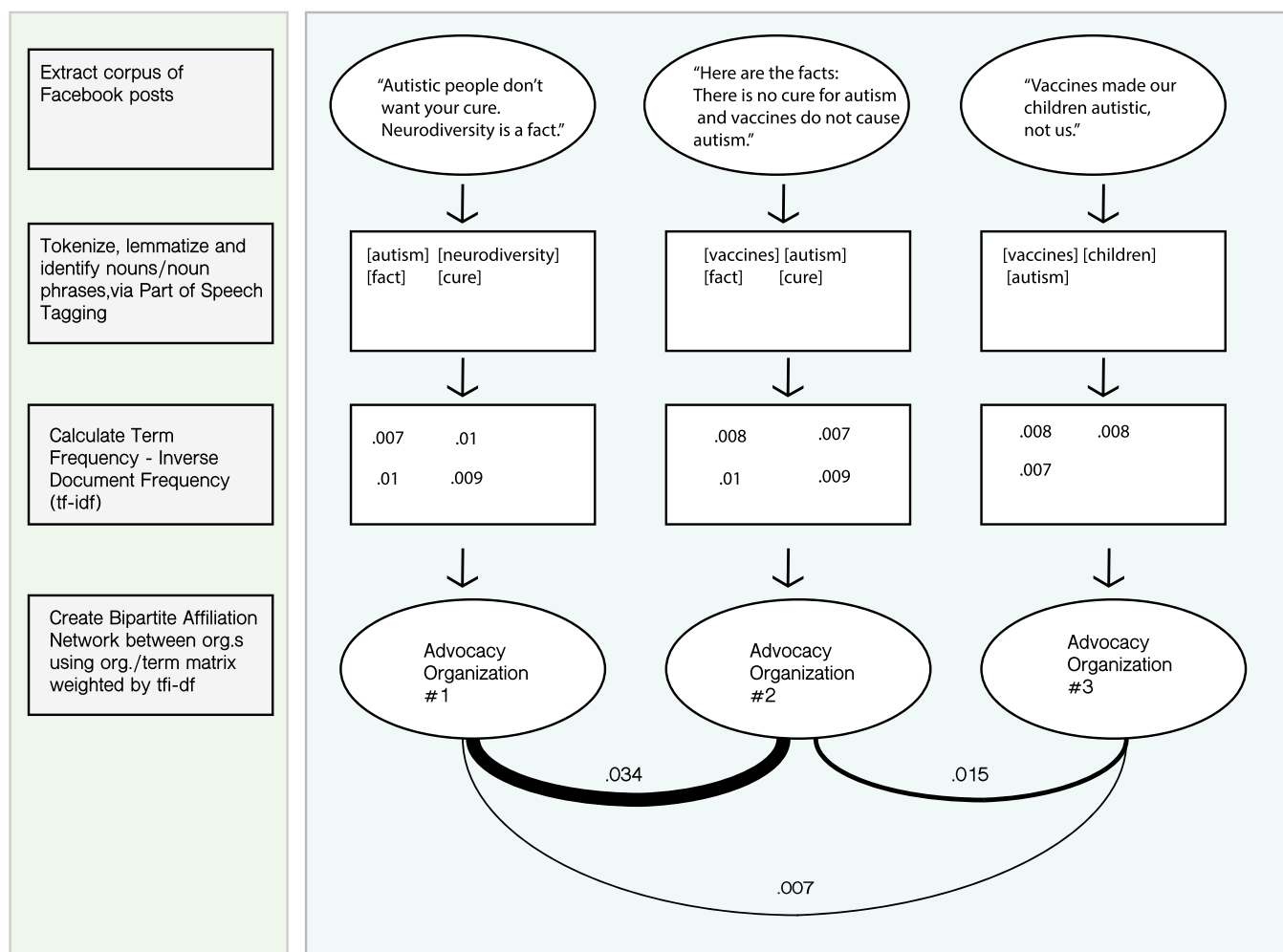
To identify organizations that create cultural bridges, I extracted all public discourse about ASDs produced by advocacy organizations and their audiences on Facebook “fan pages.” Unlike private Facebook pages, the text of all posts and comments from fan pages is publicly available because their purpose is to call public attention toward a cause.

I identified 134 Facebook pages maintained by ASD advocacy organizations, using a three-stage sampling procedure (*Materials and Methods*; see also Fig. 4). I extracted the text of all posts produced by these organizations between August 2011 and December

2012, as well as public comments about them and the names of those who made such comments, from Facebook’s Application Programming Interface. I used these data to create my outcome measure: the number of unique people who made comments about one of an organization’s posts that were more than three words who had not previously commented on any of the organization’s posts.

As Fig. 2 shows, I combined natural language processing and social network analysis to map cultural networks and identify organizations that engaged in cultural bridging during the study period. I applied standard text preprocessing techniques to the corpus of Facebook posts produced by ASD advocacy organizations. I then constructed a bipartite affiliation network linking organizations based on the amount of overlap in the nouns, proper nouns, or noun phrases they used in their Facebook posts on a given day. Finally, I used a community detection algorithm to measure the “cultural betweenness” of each organization within the cultural network, or the extent to which an organization’s message connects discursive themes that are seldom discussed together. For additional details about these procedures, see *Materials and Methods*.

Producing messages that create cultural bridges is but one of many possible ways advocacy organizations might inspire comments from new social media audiences. For example, an organization’s capacity to inspire comments may be shaped by its size, its financial resources, characteristics of its audience, or broader external factors such as the amount of public interest in ASDs on a given day. Although publicly available Facebook data are well-suited to map cultural networks and identify cultural bridges, these texts cannot be used to evaluate such alternative explanations. I therefore created a web-based Facebook application called Find Your People that offered ASD advocacy organizations a complimentary assessment of their social media strategy in return for completing an online survey about their organization and sharing nonpublic Facebook Insights data that provide more than 100 variables that describe aggregate characteristics of each organization’s audience. Institutional review boards at three universities determined this research design was exempt from human subjects review because it did not collect nonpublic information about individual social media users or manipulate them in any manner. Nevertheless, the Facebook application



**Fig. 2.** Combining natural language processing and network analysis to create cultural networks between organizations.

obtained informed consent from representatives of each organization via an authentication dialogue that directed users to a webpage that described the procedures used to protect each organization's data.

I used the data collected via the Find Your People application to create a longitudinal data set with organization per day observations. I use the following variables to control for organizational characteristics: (i) total number of Facebook fans (by day); (ii) closeness centrality of organization among all those who post, comment, or like within the social media advocacy field (by day and previous day); (iii) between-ness centrality of organization within the social media advocacy field (by day and previous day); (iv) number of page views the organization received by paying for Facebook advertising (by day); and (v) number of audiovisuals the organization used in its Facebook posts (by day). I also included the following variables to control for characteristics of Facebook audiences: (i) total number of people who viewed an organization's Facebook page (by day), (ii) percentage of page viewers younger than 35 y (by day), and (iii) percentage of page viewers from Eastern, Southern, Mid-western, and Western regions of the United States (by day). Finally, I created the following variables to measure external factors that may create opportunities for ASD advocacy organizations to stimulate public conversation: (i) number of articles about an organization in the Google News database (by day and previous day), (ii) number of blog posts about an organization listed in the Google Blogs database (by day and previous day), and (iii) relative volume of Google searches for the term "autism" (by day and previous day).

Because of significant kurtosis in the outcome variable and the uneven distribution of error across advocacy organizations, I used a negative binomial regression model with fixed effects for each organization and day. All models also controlled for the number of posts an organization produced during the previous day to ensure that the cultural bridge measure does not simply measure the effect of producing a message, regardless of its position within the cultural network.

## Results

I begin by presenting descriptive results for the entire population of ASD advocacy organizations that maintained active Facebook fan pages during the sampling period. These 134 advocacy organizations produced 28,606 messages between August 2, 2011, and December 18, 2012, which received 79,193 comments and ~1.2 million likes. The median number of daily comments more than three words in length that were produced by those who had not previously commented on an advocacy organization's posts was 0, confirming theoretical expectations that the majority of messages produced by ASD advocacy organizations receive little or no attention from new audiences. Conversely, a single organization received 883 substantial comments from new audiences in a single day.

Among the organizations in the target sample, 33.8% installed the Find Your People application in November 2012. There was very minimal evidence of response bias (see [SI Materials and Methods](#)). Fig. 3 presents the results of the fixed-effects negative binomial models. This illustration shows that the size of cultural





Finally, this study has important implications for the nascent field of computational social science. I introduced new techniques for mapping cultural networks and public deliberation that build on recent attempts to synchronize social network analysis and natural language processing (14, 15). These new tools offer a more dynamic alternative to topic models and other new forms of automated text analysis that are also poorly suited for short texts such as social media messages. Lastly, I introduced an innovative research design that used a social media application to collect a large amount of information that places social media discourse within broader social context. These methods hold considerable promise for combining the strengths of social media data with more conventional survey techniques to improve the rigor and validity of computational social science.

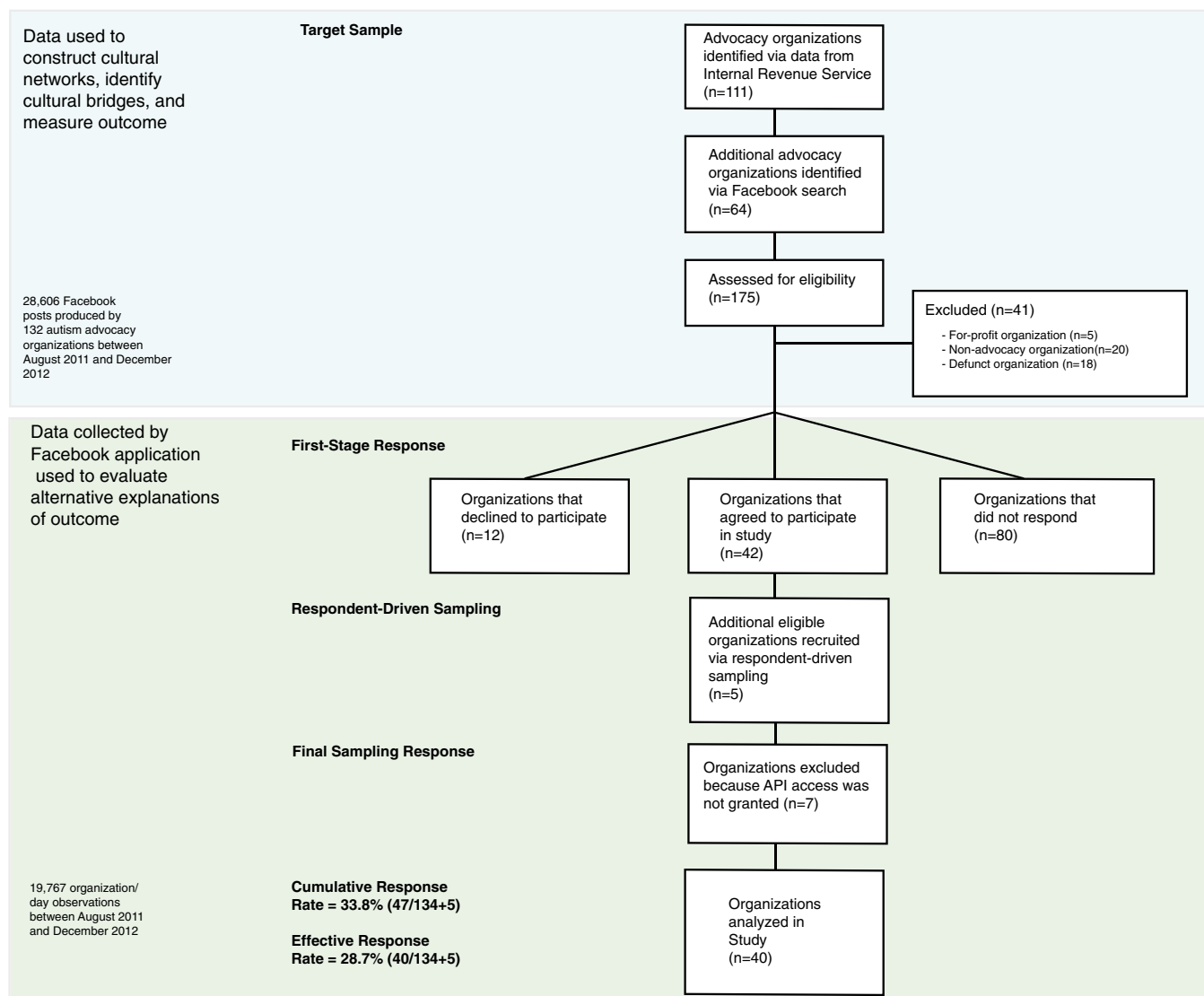
## Materials and Methods

**Sampling Advocacy Organizations.** I identified 111 autism spectrum disorder (ASD) organizations via a database of Internal Revenue Service 501(c)(3) non-profit organizations. Because advocacy organizations must have significant financial resources to register for official nonprofit status, I also used a team of research assistants to conduct searches for the terms “autism” and “Asperger’s,” using Facebook’s search function. As Fig. 4 shows, 175 organizations were identified using these two procedures, but 41 of them were excluded from the

target sample because they were not advocacy organizations or because they had recently become defunct. I used the text of posts from the remaining 134 advocacy organizations to construct daily cultural networks and identify those that created cultural bridges, using the techniques described here.

As Fig. 4 shows, 42 of the remaining 134 organizations agreed to participate in the study, and 5 additional organizations were identified via respondent-driven sampling as organizations publicized the Find Your People Facebook application with one another. The cumulative response rate was thus 33.8%. However, 7 of the 47 organizations that participated in the study submitted their data via email instead of using the Facebook application. Because significant amounts of data could not be obtained from these organizations via Facebook's Application Programming Interface, they were excluded from the study.

To identify possible sample response bias related to the study's use of a complimentary social media audit as an incentive for organizations to share their data, I examined all available data for organizations in and outside the study sample from the Internal Revenue Service and Facebook. As [Fig. S1](#) shows, no significant differences were identified according to the number of substantial comments by new social media users, the organization's annual budget, or its age. Very small, yet significant, differences were detected in the betweenness centrality measure used to generate the indicator of cultural bridges between organizations in and outside the study sample. I discuss the techniques used to account for this bias in [SI Materials and Methods](#).



**Fig. 4.** Three-stage sampling process used to recruit advocacy organizations to install Facebook application.

**Constructing Cultural Networks.** I combined four stages of natural language processing and network analysis to map cultural networks and identify advocacy organizations that created cultural bridges. Facebook posts were tokenized and then lemmatized, a process that replaces each word with its most basic syntactic form, or “lemma.” For example, the lemma of the term “running” is “run.” Next, part-of-speech taggers were applied to the lemmatized text to identify nouns, proper nouns, and noun phrases, which are most likely to define the substantive content or discursive theme of a message (15, 16). Next, I measured the term frequency–inverse document frequency for each of these terms. I then created a bipartite affiliation network that links organizations to each other on the basis of the copresence of terms within their messages (17). The weight of edges between each organization is defined by the sum of the term frequency–inverse document frequency for the overlapping terms.

Next, I calculated the betweenness centrality of each organization within this cultural network, using a variant of Dijkstra’s algorithm (18). Using this approach, the shortest path ( $d$ ) between two nodes ( $i, j$ ) can be defined as follows:

$$d^w(i, j) = \min \left( \frac{1}{w_{ih}} + \dots + \frac{1}{w_{hj}} \right) \quad [1]$$

where  $w$  is the weight of the tie between nodes and  $\alpha$  is a tuning parameter that can assume values between 0 and 1 that reflect the influence of edge weights. I used a value of 0.5, although other values of this tuning parameter produced nearly identical results. The cultural betweenness of a node ( $C$ ) during day ( $t$ ) is given by

$$C_t(i) = \frac{\sum g_{ij(i)}}{g_{ij}}, \quad [2]$$

where  $g$  is the sum of its shortest paths that pass through node  $i$  as a proportion of all shortest paths in the network.

**ACKNOWLEDGMENTS.** I thank Paul DiMaggio, James Moody, and Martin Ruef for helpful comments on the manuscript. This study was supported by the National Science Foundation and Robert Wood Johnson Foundation for financial support.

1. Developmental Disabilities Monitoring Network Surveillance Year 2010 Principal Investigators; Centers for Disease Control and Prevention (2014) Prevalence of autism spectrum disorder among children aged 8 years—Autism and Developmental Disabilities Monitoring Network, 11 sites, United States, 2010. *MMWR Surveill Summ* 63(2):1–21.
2. Bearman P (2009) Just-so stories: Vaccines, autism, and the single-bullet disorder. *Soc Psychol Q* 73(2):112–115.
3. Eyal G, Hart B, Oncular E, Oren N, Rossi N (2010) *The Autism Matrix* (Polity, Cambridge, UK).
4. Buescher AV, Cidav Z, Knapp M, Mandell DS (2014) Costs of autism spectrum disorders in the United Kingdom and the United States. *JAMA Pediatr* 168(8):721–728.
5. Dewey J (1927) *The Public and Its Problems* (Swallow Press, Athens, OH).
6. Bakshy E, Messing S, Adamic LA (2015) Political science. Exposure to ideologically diverse news and opinion on Facebook. *Science* 348(6239):1130–1132.
7. Hindman M (2009) *The Myth of Digital Democracy* (Princeton Univ Press, Princeton, NJ).
8. Sunstein CR (2002) *Republic.com* (Princeton Univ Press, Princeton, NJ).
9. White HC (1992) *Identity and Control* (Princeton Univ Press, Princeton, NJ).
10. Bearman P, Faris R, Moody J (1999) Blocking the future: New solutions for old problems in historical social science. *Soc Sci Hist* 23(4):501–533.
11. DiMaggio P (2011) Culture Networks. *The Sage Handbook of Social Network Analysis*, eds Scott J, Carrington PJ (Sage Publications, London).
12. Pachucki MA, Breiger RL (2010) Cultural holes: Beyond relationality in social networks and culture. *Annu Rev Sociol* 36(1):205–224.
13. Mohr J (1998) Measuring meaning structures. *Annu Rev Sociol* 24:345–370.
14. Vilhena DA, et al. (2014) Finding cultural holes: How structure and culture diverge in networks of scholarly communication. *Sociol Sci* 1(1):221–238.
15. Rule A, Cointet JP, Bearman PS (2015) Lexical shifts, substantive changes, and continuity in State of the Union discourse, 1790–2014. *Proc Natl Acad Sci USA* 112(35):10837–10844.
16. Borrett S, Moody J, Edelman A (2014) The rise of network ecology: Maps of topic diversity and scientific collaboration. *Ecol Model* 293:111–127.
17. Breiger RL (1974) The duality of persons and groups. *Soc Forces* 53(2):181–190.
18. Opsahl T, Agneessens F, Skvoretz J (2010) Node centrality in weighted networks: Generalizing degree and shortest paths. *Soc Networks* 32(3):245–251.
19. Allison PD, Waterman RP (2002) Fixed-effects negative binomial regression models. *Sociol Methodol* 32(1):247–265.
20. Spital A (2001) Public attitudes toward kidney donation by friends and altruistic strangers in the United States. *Transplantation* 71(8):1061.
21. 42 USC §1320b-10 (2012).
22. Healy K (2008) *Last Best Gifts Altruism and the Market for Human Blood and Organs* (Univ of Chicago Press, Chicago).
23. Simmons RG, Marine SK, Simmons RL (1987) *Gift of Life: The Effect of Organ Transplantation on Individual, Family, and Societal Dynamics* (Transaction Publishers, New Brunswick, NJ).
24. Silverman C (2013) *Understanding Autism: Parents, Doctors, and the History of a Disorder* (Princeton Univ Press, Princeton, NJ).

# Supporting Information

Bail 10.1073/pnas.1607151113

## SI Materials and Methods

**Negative Binomial Regression Models.** As Fig. S2 shows, the outcome variable was heavily skewed. No substantial comments were made by new audiences during 89.35% of the organization per day observations in the dataset.

Because of the heavy skew in the outcome variable, I examined the association between the number of comments ( $y$ ) an organization ( $i$ ) received each day ( $t$ ) and the indicators described in the main text, using a negative binomial regression model with fixed effects for time. Because fixed-effects negative binomial regression models allow for individual-specific variation in the dispersion parameter, rather than variation in the conditional mean, I also include fixed effects for each organization to facilitate unconditional estimation (19):

$$y_{it} = \alpha + \beta^T x_{it} + \mu_i + \epsilon_{it}.$$

In addition to the indicators described in the main text of this article, the model also includes an Inverse Mills Ratio variable generated using two-stage estimation procedures to account for the modest difference in the rate of cultural betweenness among organization inside and outside the final study sample. As an additional robustness check, I ran a zero-inflated regression model, a fixed-effects Poisson regression model, and a log-linear model. Each of these alternative models produced nearly identical results.

**Lag Analysis.** To identify the appropriate time interval to measure cultural betweenness, I calculated this metric at 1-, 2-, 3-, and 7-d intervals. Fig. S3 shows that the findings are highly similar for each of these periods. A 1-d interval was used in the final model because it minimized the Akaike Information Criterion, as Fig. S4 shows.

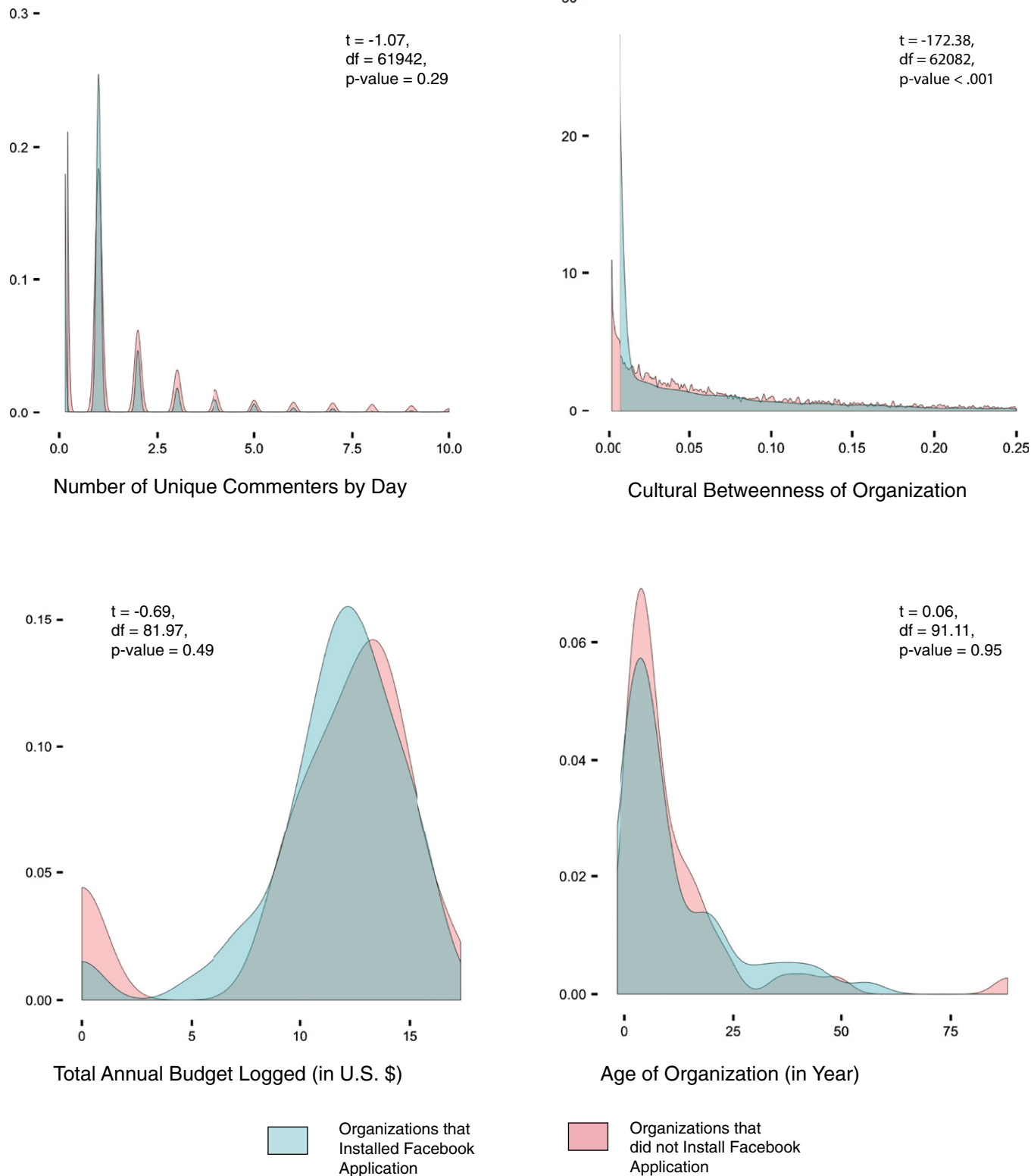
**Examining Cultural Bridges in Another Advocacy Field.** The field of ASD advocacy is highly polarized because of divergent views

about what causes ASDs and whether or how they might be cured. To determine whether the polarization of this field creates abnormal opportunities for cultural bridging, I conducted a parallel analysis of a “least similar” public health field: organ donation advocacy.

Organ donation was once viewed as a risky or unusual treatment for chronic diseases such as liver cancer or diabetes. However, organ transplants have since become the preferred treatment for these and many other chronic conditions because of recent advances in medicine. A 2001 study indicated that 92% of Americans support the process of organ donation (20), and a 2012 study revealed that 95% of Americans now support the practice (21).

As is true of most advocacy fields, public discussion of organ donation is not without any controversy. Although debates occasionally emerge about how to maximize equitable allocation of human organs, for example, there is broad-scale consensus about the urgent need to recruit new organ donors and general agreement about how this should be done (22, 23). This stands in stark contrast to the field of ASD advocacy, where there is no consensus about the root causes of ASDs or whether they can or should be cured (24).

I identified a target sample of 79 organ donation advocacy organizations in the United States and recruited them to install my application, using the same procedures I used to sample ASD advocacy organizations. Of the organ donation advocacy organizations, 59.5% installed the application, generating a panel dataset with 18,472 organization per day observations. I created the same cultural betweenness measure and other measures used in my analysis of ASD advocacy organizations. As Fig. S5 shows, the cultural betweenness measure has a strong and significant relationship with the outcome within this model, indicating that my finding about cultural bridges is not limited to the field of ASD advocacy.





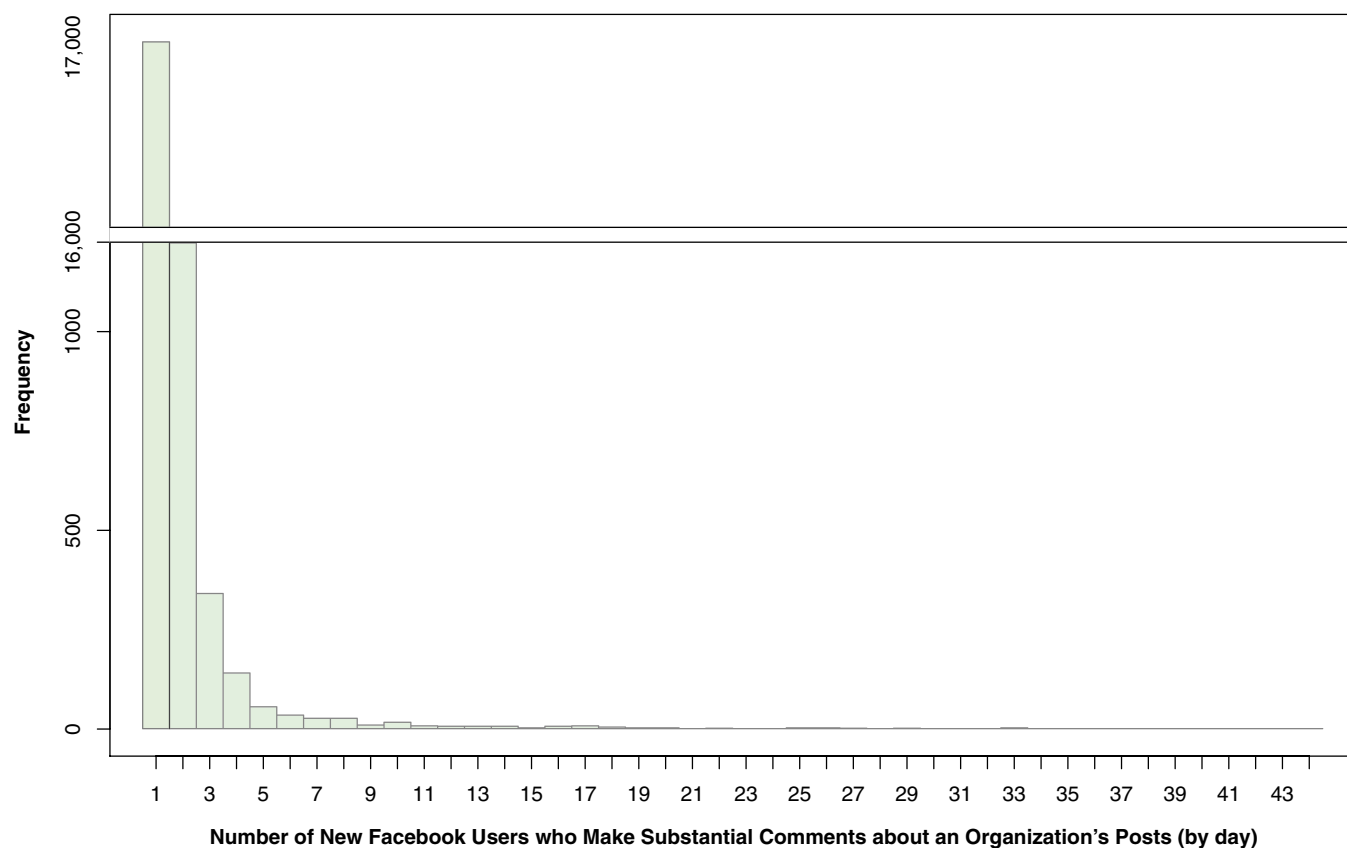


Fig. S2. Distribution of outcome variable.

