

Online Civic Intervention: A New Form of Political Participation Under Conditions of a Disruptive Online Discourse

PABLO PORTEN-CHEÉ

MARLENE KUNST

MARTIN EMMER

Freie Universität Berlin, Germany

In the everyday practice of online communication, we observe users deliberately reporting abusive content or opposing hate speech through counterspeech, while at the same time, online platforms are increasingly relying on and supporting this kind of user action to fight disruptive online behavior. We refer to this type of user engagement as *online civic intervention* (OCI) and regard it as a new form of user-based political participation in the digital sphere that contributes to an accessible and reasoned public discourse. Because OCI has received little scholarly attention thus far, this article conceptualizes low- and high-threshold types of OCI as different kinds of user responses to common disruptive online behavior such as hate speech or hostility toward the media. Against the background of participation research, we propose a theoretically grounded individual-level model that serves to explain OCI.

Keywords: political participation, disruptive online behavior, hate speech, flagging, online discussions

Current debates on revising and refocusing concepts within the field of political communication strongly consider harmful communication phenomena in online environments (Bennett & Pfetsch, 2018; Van Aelst et al., 2017). These harmful communication phenomena can take the form of an increase in so-called relativism or a felt legitimacy to question factual information (Van Aelst et al., 2017). The online world is the ecosphere where “alternative” sets of often unfounded information can be cultivated because the open architecture promotes the rapid dissemination of unverified information (Lazer et al., 2018). Another relevant harmful communication phenomenon that occurs in communicative online spaces is hate speech and incivility, which currently ranks high on the agenda of policy makers and researchers. Misogynistic and xenophobic postings especially, many of which cross the line into hate speech, are a daily occurrence (e.g.,

Pablo Porten-Cheé: p.porten-chee@fu-berlin.de

Marlene Kunst: marlene.kunst@fu-berlin.de

Martin Emmer: martin.emmer@fu-berlin.de

Date submitted: 2018-10-04

¹ This work was funded by the German Federal Ministry of Education and Research, funding code 16DII114.

Copyright © 2020 (Pablo Porten-Cheé, Marlene Kunst, and Martin Emmer). Licensed under the Creative Commons Attribution Non-commercial No Derivatives (by-nc-nd). Available at <http://ijoc.org>.

Gardiner et al., 2016). Moreover, the widespread phenomenon of trolling and the expression of fundamental hostility toward the media (Post, 2018) disrupt rational and outcome-oriented political discussions.

The comment sections of news sites have become popular spaces where users share information and opinions and deliberate on public matters. However, these often meaningful debates also attract those who are more interested in disturbing the discussion through trolling and hate speech. Yet, such utterances are in many cases sanctioned by other users who want to defend their space for discussion. In addition to interventions by activist groups such as *#ichbinhier* and *#jagärhär* (German and Swedish groups, which both translate to *#Iamhere*), individual users also react by flagging, reporting, or counterarguing uncivil, hateful, or relativistic postings. Contributing to the wider debate on updating the concepts of political communication research to identify and react to harmful communication phenomena, we introduce the concept of *online civic intervention (OCI)* and define it as a new form of user-based political participation in the digital sphere that aims to restore an accessible and reasoned public debate. On the basis of this definition, we provide an individual-level framework to identify the different conditions that promote users' engagement in OCI.

Although previous research has tried to assess why people engage in counterspeech (e.g., Kalch & Naab, 2017; Naab, Kalch, & Meitz, 2016), those studies have several shortcomings. Research on user reactions to disruptive online behavior has not yet systematized the different forms of OCI and is usually built on different research strands, hindering the development of a comprehensive OCI research field. In this study, we aim to improve our understanding of the phenomenon by first discussing some key premises of online public discourse from the theoretical perspective of deliberative democracy theory. Second, we differentiate between two OCI dimensions as new forms of political participation and propose a model that explains OCI on different levels (i.e., content, subject, context).

Online Public Discourse as a Public Matter

What can ideally be expected from online public discourse depends on the theoretical perspective. Generally, theories on democracy rely on informed citizens and a public sphere that interconnects citizens with the political decision-making process. However, the various theories on democracy deal with public communication differently. Whereas antagonistic democracy scholars argue that democracies thrive on the emotional contestation of incompatible viewpoints (see Mouffe, 1999), deliberative democracy theorists assert that discourse has to meet some standards, such as the rationality of the argumentation and comprehensiveness, to produce a thorough basis for political decision making (e.g., Fishkin, 2009; Habermas, 2006). Although all perspectives present arguable and applicable frameworks to gather the discursive realities of online discussions, we chose deliberative democracy theory as a theoretical framework because it allows for a very detailed evaluation of discussion quality. This detailed evaluation comprises a normative backbone that resonates quite well with the technological possibilities for communication online that allows literally every citizen and social group to participate and to enhance their interactions. Thus, in the following, we elaborate on the key normative demands of deliberative democracy theory for online public discourse—regardless of whether online discussions, in practice, meet the high standards of these demands.

Scholars have differentiated among the following three dimensions of deliberation: the input dimension, the discursive process dimension, and the normative output dimension (e.g., Bächtiger, Niemeyer, Neblo, Steenbergen, & Steiner, 2010). The input dimension mostly covers the design aspects in terms of providing individuals with equal opportunities to participate and express their opinions. Online discussion platforms have great potential to meet these conditions even though some societal groups are more likely than others to participate because of their digital skills and sociocultural factors (Dahlberg, 2007). With regard to the output dimension, scholars principally agree that deliberative discourse should result in a consensual solution for the issue under discussion (for an exception, see Wessler, 2008). However, because our interest is in the OCI process, both the input and outcome dimensions of deliberation theory can be neglected for now. What an improved deliberation quality will result in—for example, a less polarized antagonistic discourse, a more diverse liberal discourse, or a consensual deliberative discourse—is certainly worthy of analysis in its own right.

The process dimension of deliberation focuses on the actual conduct of discussions and is therefore the relevant guiding perspective for OCI. The measurement and comparability of the process of deliberation, however, are difficult because it is a theoretical umbrella instead of a concept (e.g., Mutz, 2008). Nevertheless, scholars have identified the key premises for deliberative discussions, which we will briefly describe.

One fundamental premise for deliberation is that each discussant has the right to question any assertion made by others, to introduce any assertion into the discourse, and to express his or her own attitudes, desires, and needs (Wessler, 2008). To ensure that everyone can make use of these rights and that no one withdraws from the discussion, the communication climate should be marked by respect and civility toward others (Wessler, 2008). Another premise for the process dimension of deliberative discourse is that opinions should be based on reasoned and rational arguments grounded on verifiable evidence or on a shared understanding of moral or normative behavior (Habermas, 2006; Stromer-Galley, 2007). Ideally, opinions should be justified by putting forward arguments with reference to external sources (Stromer-Galley, 2007). According to deliberation theory, arguments should be on-topic, as contributions from discussants should aid problem solving (Stromer-Galley, 2007). On the one hand, speakers need to consider the arguments of others when refining their own standpoint; on the other hand, they should also elaborate on their reasoning and justifications so that other discussants can incorporate the rationales into their own argumentation (Ferree, Gamson, Gerhards, & Rucht, 2002).

Many premises for public communication defined by deliberative democracy theory are intuitive parts of everyday behavior. Individuals commenting online follow norms they have appropriated during their socialization in both interpersonal interactions and media environments (de Vreese & Moeller, 2014). In other words, in situations in which users observe that other users are being attacked or intimidated by hate speech or uncivil comments, they may be motivated to intervene based on their appropriated social norms. Such social norms warrant the conditions for an accessible and reasoned public discourse at the macro level and thus refer to deliberative democratic theory.

Disrupted Conditions of Online Public Discourse

Good discourse conditions are disrupted if user contributions violate the mentioned norms. To systematize such disruptions, we examine two dimensions of online discourse conditions: namely, content and form, with the former differing between relevant and irrelevant *argumentation* and the latter between a civil and uncivil *tone*. The intersections of these two dimensions result in four online discourse conditions (Figure 1). The relevant and civil form (type I) stands for the ideal and reasoned online public discourse as discussed against the background of deliberative democracy theory. By contrast, the other three *disruptive* discourse conditions challenge an accessible and reasoned discourse. The four different types are discussed as follows.

		Content	
		Relevant	Irrelevant
Form	Civil	Ideal discourse (type I)	Irrelevant/manipulative talk (type II)
	Uncivil	Strong dispute (type III)	Hate speech (type IV)

Figure 1. Four types of (disruptive) online discourse conditions.

Civil and Irrelevant Discourse Online

Moderate harm to acceptable discourse conditions is exerted by user contributions that violate the relevance norm even though the users' expressions adhere to a civil tone (type II). Although this condition may not be directly harmful for individuals, it covers all types of manipulative argumentation and rhetoric. In this regard, today's communication science is concerned about an increase in *relativism* toward facts, evidence, and empirical knowledge (Higgins, 2016; Van Aelst et al., 2017). Assertions rooted in empirical evidence are no longer seen as facts, but instead as matters of opinion, while misinformation and conspiracy theories are on the rise (Van Aelst et al., 2017). The alternative worldviews of parts of society are supported by the dissemination of "fake news," which is spreading rapidly online (Lazer et al., 2018). Relativism poses a severe threat to argument-based discussions, as participants need to agree on facts to engage in meaningful conversations. If discussants deviate from facts and empirical knowledge, any argument becomes arbitrary.

Another critical type II case concerns *hostility against the media*. This phenomenon is based on the fact that a large number of individuals feel that they are not represented in the mass media discourse (Post, 2018). Media hostility is expressed through blanket and harsh accusations about the media, appears as part of populist communication, and can be a threat to political conversations because it is likely to distract audiences from the actual issue (Carlson, 2009).

Uncivil but Relevant Discourse Online

Another type of disruption that can call for an intervention is an uncivil but still relevant form of discussion (type III). Our notion of incivility is based on Chen's (2017) definition, which includes "insulting language or name-calling; profanity; . . . [as well as] stereotypes, and homophobic, racist, sexist, and xenophobic terms that may at times dip into hate speech" (p. 6). As this definition indicates, the line between incivility and hate speech can be blurred, but it is our understanding that incivility lacks the extreme emotion of hate speech (which falls under discourse condition IV; see also the definition in Erjavec & Kovačič, 2012). In contrast to hate speech, uncivil language does not necessarily attack individuals based on their belonging to a certain social group. In most cases, uncivil but relevant discourse online is protected by freedom-of-speech laws in democratic societies.

Uncivil and Irrelevant Discourse Online

The strongest threat to reasoned discussions online is contributions of extreme incivility that have no topical relevance, but are reduced to vilification and personal harassment (type IV). Such conditions appear in uncivil postings that, first, attack others based on their personal characteristics, and, second, because of this bias, also fail to refer to facts or other types of information that may qualify as topical contributions. According to Erjavec and Kovačič (2012), *hate speech* is an abusive, insulting, intimidating, or harassing expression that may incite violence, hatred, or discrimination based on "race, ethnic origin, religion, gender, age, physical condition, disability, sexual orientation political conviction, and so forth" (p. 900). Beyond hate speech, a special case of the uncivil and irrelevant discourse condition is *trolling*. Trolls refuse to take part in any substantive discussions and instead aim to disrupt or end discussions by silencing others (Buckels, Trapnell, & Paulhus, 2014).

To sum up, disruptive online behavior in the form of insulting, abusive, or irrational messages presents a threat to the meaningful development of a line of argumentation. Owing to such conditions, users may withdraw from online discussions or withhold their opinions for fear of being disrespectfully addressed, which might decrease the heterogeneity of participants in online discussions.

Some countries have started to regulate hate speech and fine online content providers that do not delete culpable content (Eddy & Scott, 2017). However, concerns about "overblocking" legally legitimate content have been raised (Fanta, 2018). Therefore, the quick and simple solution of eliminating problematic content runs the risk of damaging the discursive potential of the Internet (Yeo et al., 2017). Thus, a public and academic debate on collective and individual reactions to online disruptive behaviors that transcends the often ineffective or harmful approaches imposed by companies and governments is necessary.

Online Civic Intervention

Definition and Purpose

In today's online environment, the amount of user content has increased by a rate that makes it impossible to check all the content for disruptive behaviors. A part of the solution lies with the ordinary

users who are not employed by the sites, but who voluntarily take responsibility for the public discourse by engaging in OCI. We define OCI as the action taken by ordinary users to restore favorable conditions for political online discussions when threats to these conditions are perceived. OCI needs to be clearly distinguished from opinion expression. Instead, it is a form of metacommunication that calls for a more civil and reasoned public debate and aims to ensure an inclusive online public discourse. OCI helps to restore a climate of mutual respect and factual information to provide a proper space for the discussion of political matters while allowing for dissent and free speech at the same time. Finally, OCI is a highly democratic expression because the definition of what is considered a relevant and civil discourse is defined exclusively by those who are actually part of the discussion—the ordinary users—and not by a few professionals setting the platform's terms.

Through OCI, users can be involved in the process of norm-setting in online communities (Gillespie, 2018), and several aspects of OCI are a valuable complement to professional moderation. Whereas human moderators cannot realistically monitor all the discussions on social media (e.g., Facebook) in real time, users who can observe conversations and choose to intervene are always present. However, as Gillespie (2018) points out, only a small fraction of Internet users are currently likely to intervene, and this group cannot be generalized as representative of the community's will. Additionally, as we will discuss later, OCI can be misused and exploited for political ends.

Despite these concerns, we consider OCI to be an emerging and valuable form of political participation. Although political participation in a narrow sense has been traditionally defined as activities aiming to influence the actions of the government (Verba, Lehman Schlozman, & Brady, 1995), our understanding of participation changes as society changes. In various ways, the political arena has migrated online, and boundaries with other fields of civic engagement are becoming blurred (Mossberger, Tolbert, & McNeal, 2008). Communicative online spaces are where political information is frequently presented and discussed, where new forms of participation have emerged, and where political opinions are formed. One may argue that although people's intention to engage in OCI can differ (which, by the way, can apply to all other forms of participation too), OCI always contributes to an improvement in the public discourse. For example, if users choose to intervene out of commiseration with a social group being attacked and want to come to its defense, they simultaneously contribute to a nondiscriminatory communicative online space. Therefore, OCI should be considered a form of participation that enables an accessible and reasoned public discourse, which is, by its very nature, a political act.

In the following, we elaborate on the most common user behaviors related to the different forms of OCI in the literature and distinguish between *low-threshold OCI*, which is easily enabled with tools on social media platforms, and *high-threshold OCI*, which relies on a verbal and discursive form of user involvement. Low-threshold OCI not only requires less effort from users than high-threshold OCI does, but also does not expose users as much to the audience. Therefore, low-threshold OCI is assumingly more common among users than high-threshold OCI is. Some activities on platforms are designed explicitly for low-threshold OCI (e.g., flagging buttons), and others can be applied to low-threshold OCI, but were originally not designed for it (e.g., rating comments).

Low-Threshold OCI

Flagging. Flagging is probably the most common tool that platforms provide for low-threshold OCI. The opportunity to report abusive or inappropriate content through flagging is available on most social media platforms that rely heavily on user-generated content or on news sites with comment sections (e.g., Crawford & Gillespie, 2016). Although the staff makes the final decision on whether to eliminate or allow content (e.g., Crawford & Gillespie, 2016), flagging allows users to help set boundaries for what is considered acceptable speech. Therefore, it is a strategy to primarily address uncivil speech (types III and IV violations of discursive norms). From the users' perspective, flagging can be a more convenient option for reporting abusive content than replying in writing because the former is more likely to lead to immediate sanctions by professional moderators (Naab et al., 2016). Moreover, in contrast to replying to or disliking content, flagging permits users to remain anonymous. However, a disadvantage of flagging is that the sanctioned authors have no opportunity to justify themselves (Kalch & Naab, 2017). Moreover, even though some platforms allow users to give more detailed information on why they flagged specific content, users are usually not able to elaborate on the reasons for, and the context of, their concern (Gillespie, 2018). Gillespie (2018) argues that user involvement might seem more democratic than it is in reality because platforms still decide on how to handle content. Online platforms usually delete content according to whether it is punishable by law or protected by freedom-of-expression laws. Such an assessment of content is usually based on complex legal considerations, which are further complicated by internationally differing legal standards. This is the reason why, on the one hand, there may be a good reason for leaving the decision regarding content deletion up to the editors acquainted with the task. On the other hand, through flagging, users can set standards that go beyond the legal boundaries of speech and reflect the community's expectations of a healthy discourse.

Media companies clearly make use of the support they receive from their users through flagging. Twitter applies a "trusted flagger" status to users who are then given special privileges in terms of reporting abusive content. Other platforms such as YouTube and Yelp also choose to identify users who are skilled at evaluating problematic content and whose flagging is taken more seriously by professional moderators (Gillespie, 2018). Facebook applies algorithms to determine the trustworthiness of users in order to classify whose reporting behavior should be taken particularly seriously (Dwoskin, 2018).

Rating or recommending comments and users. *The New York Times* provides incentives for high-quality comments by choosing those that are particularly considerate and insightful as NYT Picks ("Comments," 2018). This task could, however, also be performed by the user community. In fact, some online news sites encourage users to rate or recommend comments or commenters (Singer, 2014). To rate other commenters, on some sites, users can even access other users' profiles to view their previous comments: For example, the platform Kialo only allows comments that include an explicit constructive argument on an issue and lets other users rate these arguments according to quality (Margolis, 2018). Consequently, commenters engaging in disruptive online behavior are affected by low ratings and can eventually be excluded from discussions. Assigning low ratings to disruptive comments can be a meaningful low-threshold OCI if users, for example, disseminate untrue information; other users will not be as likely to believe them because the low rating can make the commenter appear less credible. This procedure is inherently democratic because the community of readers and commenters decides what content is valuable

and which content can move the public discourse closer to the expectations of the audience. The collective intelligence of many users can be more suitable for determining the standards of politically meaningful conversation than the academic approach of experts, who can at times be different from the ordinary users.

Similarly, news organizations apply moderation strategies that can be adapted to promote OCI. For example, *The New York Times* "upgrades" some of its discussants to verified commenters who have the privilege of being able to comment without moderation ("Comments," 2018). Similarly, individuals who repeatedly engage in OCI can be designated as verified moderators. These individuals will get a specific recognizable profile and can caution other users if they cross the line, among others. This strategy can help increase the relative weight of constructive postings, resulting in more civil and relevant user contributions (type I). Some news sites have already embedded particular users in their moderation team (Reich, 2011).

In general, formally attributing more responsibility to certain users can promote their engagement in OCI. For example, research on the community management of list servers has shown that individuals with a formal leadership role make substantially more effort to moderate communication through list servers than ordinary members do (Butler, Sproull, Kiesler, & Kraut, 2007). However, to officially give certain users the task of moderating without compensating them financially is ethically questionable. In the case of Reddit, voluntary moderators went on strike to express their discontent with not being paid for their labor (Matias, 2019).

High-Threshold OCI

High-threshold OCIs as discursive strategies are seldom universal, but they are specific to the respective disruptive online behavior. Therefore, high-threshold and verbal OCIs must be addressed separately for the respective disruptive online behaviors.

With regard to hate speech and incivility, a good example of a high-threshold OCI is the activist community associated with the German hashtag *#ichbinhier*. This group attempts to deescalate incivility and hate speech in discussions through the use of friendly counterspeech (Ley, 2018). However, even in this group of activists, there has been some disagreement about the tolerance level and tone of the counterspeech in response to xenophobic comments, illustrating that both deviance and adequate response strategies can be judged differently. Similarly, organized collective responses to harassment have been observed in the gamer community, where women, particularly women of color, are frequently harassed (Gray, 2013). These women reacted by annoying the male harassers through destroying the game or sitting in on games and not playing. Moreover, activists in Australia have collectively drawn attention to online violence against women on social media, resulting in legal consequences for one of the perpetrators (Jane, 2017). However, engaging in such groups requires not only time and commitment, but also a great deal of civil courage, as members can be intimidated and threatened by the individuals they counter (Ley, 2018; Matias, 2019). This situation may be one of the reasons that participants are more likely to respond to hateful comments by flagging them than by responding verbally to them (Kalch & Naab, 2017).

A promising strategy against trolling, which is classified as a high-threshold OCI, is to debunk trolls by disclosing their identity and to advise other users to ignore them (Turner, Smith, Fisher, & Welsler, 2005).

A low-threshold OCI, such as disliking a comment or flagging it, may not be as efficient because trolls can rarely be identified from a single post. Warning others “not to feed the trolls” has actually been a common behavior in online discussions since the early days of the Internet (Binns, 2012).

In the case of media hostility, a fine line exists between the legitimate criticism of journalists and the blanket hostility toward mainstream media. Defending the media each time they are criticized is not reasonable because the audience is supposed to hold journalists accountable. However, when media criticism becomes all-encompassing, it redirects the discussion from the current issues. In these cases, OCI is particularly important because users may be more effective than journalists when responding to the harsh criticisms against the media; journalists may not be perceived as neutral, but as biased and defending themselves. By contrast, users can encourage other users to check multiple sources to verify whether the information provided by the news is correct. Moreover, users can clarify for others when a news item is marked as an opinion piece and is therefore being unjustly accused of lacking objectivity. If a news item is considered “fake news,” as Donald Trump and some of his supporters frequently claim, users can post URLs from a variety of other sources that either confirm or deny the credibility of the attacked item.

Responding with OCI to relativism or to the dissemination of misinformation requires sufficient knowledge about a topic. No community managers, not even those with an excellent education, can be an expert in all topics that the public is eager to discuss. Consequently, mass media and social media platforms can identify experts in their communities and involve them in fighting relativism. Some platforms already collaborate with fact-checking services and journalistic organizations with regard to misinformation, acknowledging that the review of information should be undertaken by experts (Caplan, Hanson, & Donovan, 2018; Dwoskin, 2018). In the comment sections, particularly knowledgeable users from different fields can obtain expert profiles, which are identified by other users as an indication that they are especially trustworthy sources. To some extent, the expert approach contradicts the egalitarian ideal of online discussions, but it may be necessary to maintain the proper conditions for rational discourse. Moreover, ordinary users can also provide links to trustworthy sources that counter misinformation or point out when a comment clearly contains false information. However, in this case, it is the word of one user against another.

Figure 2 summarizes the possible high- and low-threshold OCIs for a particular disruptive online behavior.

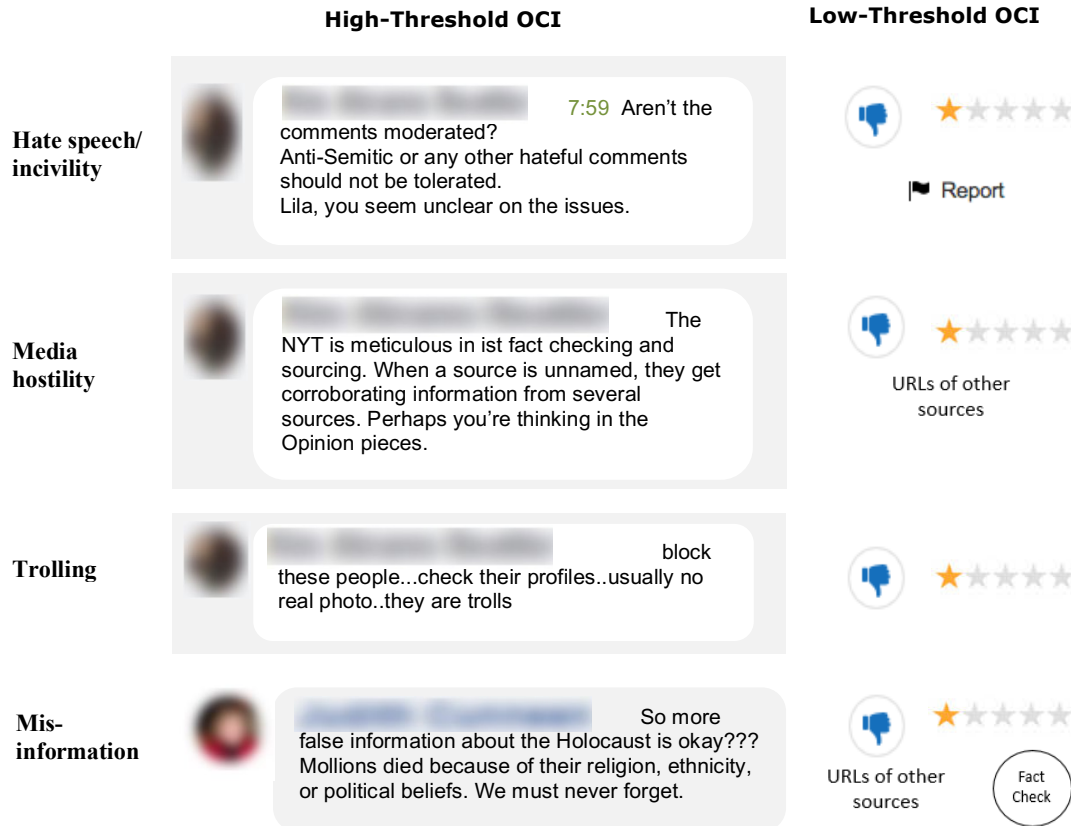


Figure 2. Examples of exemplary high- and low-threshold OCIs for different disruptive online behaviors. Comments are taken from the Facebook page of The New York Times.

Explanatory Model for OCI

We assume that individuals must cross several thresholds before they engage in OCI. First, users need to perceive online content as a violation of social norms. Although the “(potential) deviance of online comments becomes a function of [the] user’s inherent belief system” (Wilhelm & Joeckel, 2019, p. 382), some characteristics of disruptive online behaviors appear as particularly threatening, such as incitement to violence against a social group or strongly abusive language (Kalch & Naab, 2017; Leonhard, Rueß, Obermaier, & Reinemann, 2018; Wilhelm, Ziegler, & Joeckel, 2019). Second, we assume that several individual factors predict OCI. Among these are predictors that empirical research has found to promote political participation or opinion expression, but also individuals’ attitudes toward issues or their social identification with a verbally attacked group. Third, there are contextual factors in communicative online spaces that make OCI more likely. In this context, scholars have, for example, drawn on the theory of bystander behavior to analyze whether the actions of other users influence user engagement in OCI (Leonhard et al., 2018; Naab et al., 2016). In sum, we expect that the more all these requirements are

met, the more users will engage in OCI, although the factors may be interrelated. For example, if an individual possesses characteristics (individual factors) that are favorable to OCI, this individual may engage in OCI even if the perceived threat of the online content is relatively low.

In the following, we discuss each factor in detail to provide a comprehensive explanatory model for OCI that can inspire further empirical research (Figure 3). The model draws on the findings from previous studies on OCI, but also integrates related research that, for example, addresses responses to hate speech in offline settings.

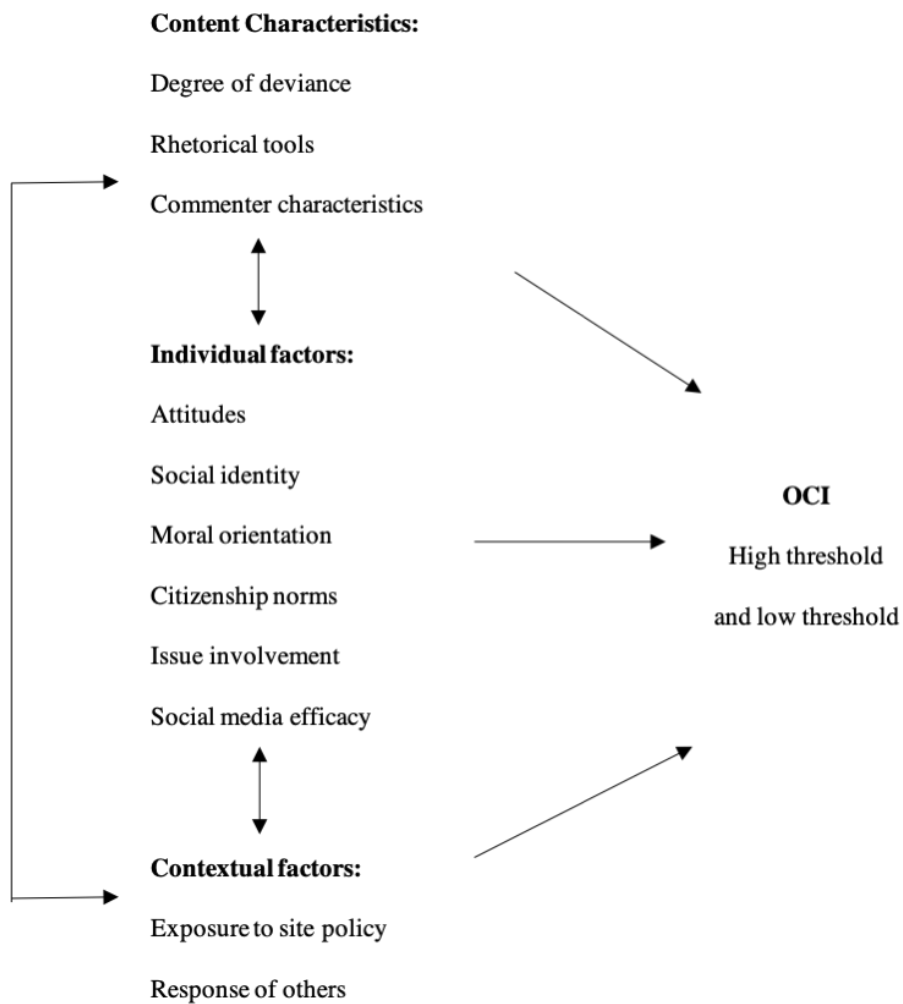


Figure 3. Assumed predictors of OCI.

Disruptive Online Behavior as a Perceived Threat

Theoretical considerations. We assume that among the three types of disruptive online behavior, uncivil content (types III and IV) is likely to be perceived as more threatening than irrelevant but civil disruptions (type II). For the latter to be perceived as threatening the discourse, users must have sufficient knowledge to identify and counter false information, which requires expert knowledge (Caplan et al., 2018). Ideally, such users also have the ability to differentiate between legitimate challenges to dominant narratives and false information, given that challenging the dominant discourse is desirable from the perspective of public sphere theory. However, the line between both can be blurred, and some users may respond to any content that challenges the status quo with OCI because they perceive that content as disruptive. Similarly, in the case of media hostility (type II), users need to differentiate between legitimate media criticism on the one hand, and sweeping media criticism that fundamentally attempts to destroy the media's credibility on the other. Whether media hostility is perceived as a threat to public discourse is likely to depend on users' personal attitudes toward the media. It will probably increase when it is accompanied by strong language and crosses the line into uncivil disruptive online behavior (type III or IV). In comparison, when trolling is civil in tone, it can rarely be identified as disruptive behavior based on a single comment. In fact, when considered individually, most of the comments made by trolls cannot be classified as trolling (Samory & Peserico, 2017; Turner et al., 2005). Therefore, trolls are less suspicious toward other readers because they "hide among neighboring posts, and build upon an existing state of excitement in the discussion, to amplify controversy" (Samory & Peserico, 2017, p. 6946). Experienced users are more likely to detect trolling and perceive it as a threat (Shaw, 2013).

With regard to the types of disruptive online behavior that include incivility or hate speech (types III and IV), three content characteristics are likely to influence whether users perceive it as a threat: the language, the social group that is verbally attacked, and the commenters' cues. In general, we assume that the more severe and the stronger the language is, the more it will be perceived as threatening, with hate speech scoring higher in both severity and strength of language than incivility. In this regard, the specific rhetorical means applied by the commenter—for example, particularly threatening language—can intensify the perceived deviance of the comment. Concerning the verbally attacked group, its perceived vulnerability tends to increase the perceived harmfulness of hate speech in offline settings (Cowan & Mettrick, 2002). Consequently, users may be more likely to respond with OCI to comments that include hate or uncivil remarks about refugees, as compared with such remarks about less vulnerable social groups, such as environmental activists or vegans. The profile of the commenter can increase the perceived deviance of the content—for example, if a comment that rages against refugees is written by a person with a pseudonym or a profile picture that clearly identifies him or her as a neo-Nazi.

Empirical findings. Experimental inquiries into comment sections online have shown that only strong and overtly offensive language is perceived as sufficiently deviant for users to engage in OCI, whereas more subtle verbal attacks on social groups are not. In their experiment, Kalch and Naab (2017) drew on Papacharissi's (2004) distinction between impoliteness and incivility, with the impolite comments including insults and vulgarity toward Muslims, and the uncivil comments stereotyping Muslims in a negative way without using abusive language. Whereas impolite comments affected flagging and counterspeech, uncivil comments did not (Kalch & Naab, 2017). As the authors noted, it was alarming that the participants did not

perceive uncivil comments as sufficiently problematic even though these comments included harsh attacks on Muslims. Similarly, studies have shown that OCI is more likely when comments openly call for violence than when hate is more subtle (Leonhard et al., 2018; Wilhelm et al., 2019). Moreover, simple rhetorical tools applied by hate commenters seem to efficiently decrease the likelihood of OCI: For example, flagging is less likely when a hate comment is accompanied by a justification made by the commenter (e.g., claiming no intent to offend others; Wilhelm et al., 2019).

Generally, whether a hate comment addresses an abstract social group or an individual appears to be irrelevant for users' engagement in OCI. However, in combination with other comment characteristics, a study by Naab et al. (2016) shows that flagging is less likely for verbal attacks on abstract groups than on individuals (see also Wilhelm et al., 2019). Whether the perceived vulnerability of a social group influences users' OCI has not been addressed yet. Similarly, empirical findings on the effect of the commenter's characteristics are rare. One exception is the study by Wilhelm and Joeckel (2019), who find that users flag hate speech by female commenters more than they do hate speech by male commenters. The authors argue that because of gender stereotypes, hate comments by women are perceived as particularly deviant because such behavior is especially unexpected from females.

Individual Factors

Theoretical considerations. Although OCI is not an opinion expression, individuals' attitudes are likely to influence whether users choose to intervene, given that attitudes affect the perceived deviance of content. For example, verbal attacks against homeless people may be perceived as less deviant if one considers homelessness as the result of individual failure rather than structural circumstances. Positive attitudes toward social groups such as homosexuals can be expected to increase the perceived deviance of verbal attacks against this group, and vice versa. A potential predictor is not only attitudes toward the social group, but also individuals' social identity: If individuals identify with a group that is verbally attacked, one can assume that they will be more sensitive to hate speech or incivility and more likely to react.

Moreover, information is more likely to be perceived as untrue or deviant if it counters one's own beliefs, and even more so if issue involvement is high. This assumption is based on the theory of cognitive dissonance, which implies that individuals experience discomfort when exposed to counterattitudinal content (Festinger, 1957). Therefore, a strategy to avoid the discomfort is to mark the counterattitudinal content as untrue. In general, the role of attitudes with regard to OCI can become problematic. For example, flagging has been misused for political means by users who systematically flag content because they disagree with it, or they want to silence or harass certain social groups (Dwoskin, 2018; Gillespie, 2018). Moreover, sensitivity toward media hostility is likely to depend strongly on the individual's general attitudes toward the media. Individuals with strong media trust will presumably perceive attacks on the media as very threatening, whereas for individuals with moderate or low media trust, the threshold for unacceptable media hostility will not be crossed easily.

Empirical findings. Concerning attitudes, the study by Kalch and Naab (2017) supports that individuals with positive attitudes toward Muslims are more likely to flag Islamophobic comments than are individuals with less positive attitudes toward them. However, with regard to social identity, the arousal of

negative emotions when being exposed to verbal attacks against a social group tends to be just as strong for in-group as for out-group members (Kang, Thorson, Fung, & Borah, 2009). Kang et al. (2009) assume that in-group members of a group that attacks the out-group experience negative emotions because they do not approve of the uncivil behavior and are disappointed by their associates (Kang et al., 2009). In this vein, Boeckmann and Liew (2002) find that the negative emotions of individuals with a strong Asian American identity are equally vigorous when exposed to hate speech against Asian Americans as against African Americans because, in the latter case, the social identity of the participants becomes more abstract as simply belonging to a minority group. In sum, although these studies do not address OCI, OCI may be driven by attacks against either one's own group or other groups with whom one presumably shares certain commonalities (e.g., minority status).

Another individual characteristic that OCI research has addressed is users' attitudes toward their moral foundations. Studies have found that an individualizing moral foundation, which focuses on the rights and protection of individuals, fosters reporting behavior, whereas a binding moral foundation, which, among other things, focuses on in-group loyalty, has been shown either to reduce it or to not have an effect (Wilhelm & Joeckel, 2019; Wilhelm et al., 2019). Thus, as individualizing moral foundations are connected with a stronger commitment to fairness and care, these individuals show stronger engagement in OCI (Wilhelm et al., 2019). However, whether individuals with an individualizing moral foundation perceive hate comments as more deviant remains unclear. Moreover, Watson, Peng, and Lewis (2019) find that young and wealthy White males are most likely to report deviant content in general. Moreover, the authors find that authoritarianism is a positive predictor of reporting abusive content, which is a way for authoritarian personalities to maintain the moral order (Watson et al., 2019). Individuals with higher trust in the media tend to flag more; only individuals with trust in the media as an authority tend to be willing to cooperate with them (Watson et al., 2019).

Potential factors. Although not yet tested, several factors are likely to promote individuals' engagement in OCI because they are general predictors of participation. For example, having an extrovert personality is a predictor of interaction in social media environments (Seidman, 2013). Moreover, political social media efficacy predicts social media activism, and this is why we argue that such a belief in one's own capability is also a likely predictor of OCI (Velasquez & Larose, 2015).

Moreover, we assume OCI to be more common among individuals with strong support for citizenship norms. Citizenship norms are a set of expectations of what individuals think people should do if they are good citizens (Dalton, 2008). Thus far, empirical studies have shown that such norms predict individuals' political participation (e.g., Dalton, 2008). Because we consider OCI to be a new type of political participation, we assume that strong support for citizenship norms is associated with OCI. However, because OCI is not a conventional type of political participation, we expect that conventional types of citizenship norms (e.g., the expectation that "one should vote in every election") will have a weak effect. By contrast, the expectation that a "good citizen" should engage in political discussions is likely to be a predictor for OCI.

Research on whether such individual characteristics predict engagement in OCI is rare, and thus, this calls for more scholarly attention.

Contextual Factors

Theoretical considerations. Last, some contextual factors can promote OCI. One of these is the available instruments that users can apply to report abusive or untrue content. Other contextual factors likely to have an impact are the community norms negotiated in group-based online environments, which are observed and appropriated through constant observation of what other peers comment on and how. In this vein, Shaw (2013) observes that "new people coming into communities develop an awareness of their right to disallow harassment and offensive comments in their own blogs by observing moderation practices on other blogs" (p. 106). Likewise, community norms can promote OCI when an individual sees others doing it. Conversely, that other users may potentially intervene when exposed to disruptive online behavior can also reduce one's own perceived responsibility to engage in OCI.

Empirical findings. Flagging tends to increase when individuals are exposed to information on usage policies, intervention options, and explicit encouragement to report norm violations (Naab et al., 2016). Therefore, the explicit encouragement of users to intervene appears to be an effective strategy for online sites to involve more users in moderation.

Studying how other users' interventions affect one's own, Naab et al. (2016) found that users' intention to flag disruptive comments increases when others have not previously sanctioned such comments. But as soon as users perceive others to have reacted to disruptive comments with OCI, their perceived self-responsibility to engage in OCI decreases. Naab et al. (2016) assume that other users' OCI leads to the impression that someone else is taking care of the situation. In addition, the perceived reach of the disruptive comments also influences OCI intention. In this vein, participants of Leonhard and colleagues' (2018) study who saw that many other users had been exposed to a hate comment had a lower intention to intervene than those who saw that only a few other users had witnessed the incident. Thus, the studies by Naab et al. (2016) and Leonhard et al. (2018) both find support for the bystander theory.

Conclusion

Although journalists are increasingly taking on the role of community managers (e.g., Bakker, 2014), media companies outsource moderation to professional content moderators (Ihlebaek & Krumsvik, 2015) because taking action against disruptive online behavior is an enormous task. Therefore, media platforms have begun to acknowledge that they need the assistance of ordinary users for moderation. Thus, many encourage low-threshold OCI—for example, by providing buttons to flag abusive comments. Gillespie (2018) criticizes the involvement of ordinary users as unpaid labor because it is unethical and unreasonable. However, letting users participate in defining and enforcing the social norms of online discourse is also inherently democratic. Encouraging such behavior can be highly beneficial for the online public sphere that sets its normative standards beyond the boundaries imposed by governments and corporations.

In this study, we proposed to distinguish between low- and high-threshold OCIs to provide a more systematized approach to studying OCI. However, the lines between both forms can be blurred. For instance, when exposed to expressed media hostility, individuals may first encourage others to check

multiple sources (high-threshold OCI) and subsequently post URLs that lead others directly to such sources (low-threshold OCI).

Through social learning processes, OCI can establish itself as a new form of political participation. Our perspective is certainly optimistic, and we are well aware that each type of OCI can be misused and lead to even more disruptions in the public sphere. Hate comments can be rated as valuable by like-minded users, and users can post URLs from fake media sites that supposedly provide evidence for biased coverage by the mainstream media. Moreover, not every user will agree with those who insist on a respectful tone in the discussion. For example, right-wing nationalists are likely to discredit users who fight against hate speech targeting ethnic minorities as “social justice warriors” and to counteract their efforts (Lewis, 2018). Therefore, individuals who engage in OCI need to be persistent and able to withstand resistance. However, we believe that the atmosphere in many online discussions will shift if more people view it as their civic duty to engage in OCI. OCI should receive more attention as part of digital literacy education. One such example is the No Hate Speech Youth Campaign of the Council of Europe that encourages young people to engage in OCI.

Understanding how OCI works can be achieved using a wide variety of methodological approaches. Through qualitative and quantitative content analyses of online discussions, researchers can analyze when and how users engage in high-threshold OCI. Moreover, digital trace data—including flagging, comment ratings, and fact-checking requests—can produce insights into which comment characteristics trigger which specific type of low-threshold OCI. Using surveys or experiments, researchers could explore when and which individuals apply the available instruments with an OCI intention. Scholars should be encouraged to include items that reflect low- and high-threshold OCIs in survey research on political participation. Another question that both survey and experimental research could address is what makes individuals predisposed to using particular types of OCI. We do, for instance, assume that individuals who strongly follow citizenship norms (e.g., Dalton, 2008) may be particularly likely to engage in this type of political participation. Finally, psychological characteristics such as conflict avoidance can make individuals more likely to opt for low-threshold OCI than for high-threshold OCI or less likely to engage in OCI altogether. Although some scholars have already started to address this body of questions, more clarifying inquiries into this underresearched but highly relevant field are necessary to find ways to deal with disruptions in online public discourse.

References

- Bächtiger, A., Niemeyer, S., Neblo, M., Steenbergen, M. R., & Steiner, J. (2010). Disentangling diversity in deliberative democracy: Competing theories, their blind spots and complementarities. *Journal of Political Philosophy*, 18(1), 32–63. doi:10.1111/j.1467-9760.2009.00342.x
- Bakker, P. (2014). Mr. Gates returns: Curation, community management and other new roles for journalists. *Journalism Studies*, 15(5), 596–606. doi:10.1080/1461670X.2014.901783

- Bennett, W. L., & Pfetsch, B. (2018). Rethinking political communication in a time of disrupted public spheres. *Journal of Communication*, 68(2), 243–253. doi:10.1093/joc/jqx017
- Binns, A. (2012). Don't feed the trolls! Managing troublemakers in magazines' online communities. *Journalism Practice*, 6(4), 547–562. doi:10.1080/17512786.2011.648988
- Boeckmann, R. J., & Liew, J. (2002). Hate speech: Asian American students' justice judgments and psychological responses. *Journal of Social Issues*, 58(2), 363–381.
- Buckels, E. E., Trapnell, P. D., & Paulhus, D. L. (2014). Trolls just want to have fun. *Personality and Individual Differences*, 67, 97–102. doi:10.1016/j.paid.2014.01.016
- Butler, B., Sproull, L., Kiesler, S., & Kraut, R. E. (2007). Community effort in online groups: Who does the work and why? In S. Weisband (Ed.), *Leadership at a distance: Research in technologically supported work* (pp. 346–362). Hillsdale, NJ: Erlbaum.
- Caplan, R., Hanson, L., & Donovan, J. (2018). Dead reckoning: Navigating content moderation after "fake news." *Data & Society*. Retrieved from https://datasociety.net/pubs/oh/DataAndSociety_Dead_Reckoning_2018.pdf
- Carlson, M. (2009). Media criticism as competitive discourse: Defining reportage of the Abu Ghraib scandal. *Journal of Communication Inquiry*, 33(3), 258–277. doi:10.1177/0196859909333693
- Chen, G. M. (2017). *Nasty talk: Online incivility and public debate*. Cham, Switzerland: Palgrave Macmillan.
- Comments. (2018). *The New York Times*. Retrieved from <https://help.nytimes.com/hc/en-us/articles/115014792387-Comments>
- Cowan, G., & Mettrick, J. (2002). The effects of target variables and setting on perceptions of hate speech. *Journal of Applied Social Psychology*, 32(2), 277–299. doi:10.1111/j.1559-1816.2002.tb00213.x
- Crawford, K., & Gillespie, T. (2016). What is a flag for? Social media reporting tools and the vocabulary of complaint. *New Media & Society*, 18(3), 410–428. doi:10.1177/1461444814543163
- Dahlberg, L. (2007). The Internet, deliberative democracy, and power: Radicalizing the public sphere. *International Journal of Media and Cultural Politics*, 3(1), 47–64. doi:10.1386/macp.3.1.47/1
- Dalton, R. J. (2008). Citizenship norms and the expansion of political participation. *Political Studies*, 56(1), 76–98. doi:10.1111/j.1467-9248.2007.00718.x

- de Vreese, C. H., & Moeller, J. (2014). Communication and political socialization. In C. Reinemann (Ed.), *Political communication* (pp. 529–546). Berlin, Germany: Walter de Gruyter.
- Dwoskin, E. (2018, August 21). Facebook is rating the trustworthiness of its users on a scale from zero to 1. *Washington Post*. Retrieved from <https://www.washingtonpost.com/technology/2018/08/21/facebook-is-rating-trustworthiness-its-users-scale-zero-one/>
- Eddy, M., & Scott, M. (2017, June 30). Delete hate speech or pay up, Germany tells social media companies. *The New York Times*. Retrieved from <https://www.nytimes.com/2017/06/30/business/germany-facebook-google-twitter.html>
- Erjavec, K., & Kovačič, M. P. (2012). "You don't understand, this is a new war!": Analysis of hate speech in news web sites' comments. *Mass Communication & Society*, 15(6), 899–920.
- Fanta, A. (2018, May 3). EU-Parlament warnt vor Overblocking durch Internetfirmen [EU-Parliament warns of overblocking by Internet companies]. *Netzpolitik*. Retrieved from <https://netzpolitik.org/2018/eu-parlament-warnt-vor-overblocking-durch-internetfirmen/>
- Ferree, M. M., Gamson, W. A., Gerhards, J., & Rucht, D. (2002). Four models of the public sphere in modern democracies. *Theory and Society*, 31(3), 289–324. doi:10.1023/A:1016284431021
- Festinger, L. (1957). *A theory of cognitive dissonance*. Stanford, CA: Stanford University Press.
- Fishkin, J. S. (2009). *When the people speak: Deliberative democracy and public consultation*. Oxford, UK: Oxford University Press.
- Gardiner, B., Mansfield, M., Anderson, I., Holder, J., Louter, D., & Ulmanu, M. (2016, April 12). The dark side of Guardian comments: The Web we want. *The Guardian Online*. Retrieved from <https://www.theguardian.com/technology/2016/apr/12/the-dark-side-of-guardian-comments>
- Gillespie, T. (2018). *Custodians of the Internet: Platforms, content moderation, and the hidden decisions that shape social media*. New Haven, CT: Yale University Press.
- Gray, K. L. (2013). Collective organizing, individual resistance, or asshole grievers? An ethnographic analysis of women of color in Xbox Live. *Ada: A Journal of Gender, New Media, and Technology*, 2. Retrieved from <https://adanewmedia.org/2013/06/issue2-gray/>
- Habermas, J. (2006). Political communication in media society: Does democracy still enjoy an epistemic dimension? The impact of normative theory on empirical research. *Communication Theory*, 16(4), 411–426. doi:10.1111/j.1468-2885.2006.00280.x
- Higgins, K. (2016, November 28). Post-truth: A guide for the perplexed. *Nature*, 540(9). doi:10.1038/540009a

- Ihlebaek, K. A., & Krumsvik, A. H. (2015). Editorial power and public participation in online newspapers. *Journalism*, 16(4), 470–487. doi:10.1177/1464884913520200
- Jane, E. A. (2017). Feminist digilante responses to a slut-shaming on Facebook. *Social Media+ Society* 3(2). doi:10.1177/2056305117705996
- Kalch, A., & Naab, T. K. (2017). Replying, disliking, flagging: How users engage with uncivil and impolite comments on news sites. *Studies in Communication and Media*, 6(4), 395–419. doi:10.5771/2192-4007-2017-4-395
- Kang, N., Thorson, K., Fung, T. K. F., & Borah, P. (2009, May). *Uncivil engagement: Linking incivility to political participation through negative emotions*. Paper presented at the annual conference of the International Communication Association, Chicago, IL. Retrieved from http://citation.allacademic.com/meta/p_mla_apa_research_citation/3/0/0/9/9/pages300993/p300993-1.php
- Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., . . . Zittrain, J. L. (2018). The science of fake news. *Science*, 359(6380), 1094–1096. doi:10.1126/science.aao2998
- Leonhard, L., Rueß, C., Obermaier, M., & Reinemann, C. (2018). Perceiving threat and feeling responsible. How severity of hate speech, number of bystanders, and prior reactions of others affect bystanders' intention to counterargue against hate speech on Facebook. *Studies in Communication and Media*, 7(4), 555–579. doi:10.5771/2192-4007-2018-4-555
- Lewis, R. (2018). *Alternative influence: Broadcasting the reactionary right on YouTube*. Retrieved from <https://datasociety.net/output/alternative-influence/>
- Ley, H. (2018). *#ICHBINHIER. Zusammen gegen fake news und Hass im Netz [#IAMHERE. Together against fake news and hate online]*. Cologne, Germany: DuMont Buchverlag.
- Margolis, J. (2018, January 24). Meet the start-up that wants to sell you civilised debate. *Financial Times*. Retrieved from <https://www.ft.com/content/4c19005c-ff5f-11e7-9e12-af73e8db3c71>
- Matias, J. N. (2019). The civic labor of volunteer moderators online. *Social Media + Society*, 5(2). doi:10.1177/2056305119836778
- Mossberger, K., Tolbert, C. J., & McNeal, R. S. (2008). *Digital citizenship: The Internet, society and participation*. Cambridge, MA: MIT Press.
- Mouffe, C. (1999). Deliberative democracy or agonistic pluralism? *Social Research*, 66(3), 745–758.
- Mutz, D. C. (2008). Is deliberative democracy a falsifiable theory? *Annual Review of Political Science*, 11, 521–538. doi:10.1146/annurev.polisci.11.081306.070308

- Naab, T. K., Kalch, A., & Meitz, T. G. K. (2016). Flagging uncivil user comments: Effects of intervention information, type of victim, and response comments on bystander behavior. *New Media and Society, 20*(2), 777–795. doi:10.1177/1461444816670923
- Papacharissi, Z. (2004). Democracy online: Civility, politeness, and the democratic potential of online political discussion groups. *New Media & Society, 6*, 259–283. doi:10.1177/1461444804041444
- Post, S. (2018). Polarizing communication as media effects on antagonists: Understanding communication in conflicts in digital media societies. *Communication Theory, 29*(2), 213–235. doi:10.1093/ct/qty022
- Reich, Z. (2011). User comments. In J. B. Singer, A. Hermida, D. Domingo, A. Heinonen, S. Paulussen, T. Quandt . . . M. Vujnovic (Eds.), *Participatory journalism* (pp. 96–117). Malden, MA: Wiley-Blackwell.
- Samory, M., & Peserico, E. (2017). Sizing up the troll. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (pp. 6943–6947). New York, NY: ACM. doi:10.1145/3025453.3026007
- Seidman, G. (2013). Self-presentation and belonging on Facebook: How personality influences social media use and motivations. *Personality and Individual Differences, 54*(3), 402–407. doi:10.1016/j.paid.2012.10.009
- Shaw, F. (2013). Still “searching for safety online”: Collective strategies and discursive resistance to trolling and harassment in a feminist network. *The Fibreculture Journal, 22*, 93–108. Retrieved from <http://twentytwo.fibreculturejournal.org/fcj-157-still-searching-for-safety-online-collective-strategies-and-discursive-resistance-to-trolling-and-harassment-in-a-feminist-network/>
- Singer, J. B. (2014). User-generated visibility: Secondary gatekeeping in a shared media space. *New Media & Society, 16*(1), 55–73. doi:10.1177/1461444813477833
- Stromer-Galley, J. (2007). Measuring deliberation’s content: A coding scheme. *Journal of Public Deliberation, 3*(1), 1–35.
- Turner, T. C., Smith, M. A., Fisher, D., & Welsler, H. T. (2005). Picturing Usenet: Mapping computer-mediated collective action. *Journal of Computer-Mediated Communication, 10*(4). doi:10.1111/j.1083-6101.2005.tb00270.x
- Van Aelst, P., Strömbäck, J., Aalberg, T., Esser, F., de Vreese, C., Matthes, J., . . . Stanyer, J. (2017). Political communication in a high-choice media environment: A challenge for democracy? *Annals of the International Communication Association, 41*(1), 3–27. doi:10.1080/23808985.2017.1288551

- Velasquez, A., & LaRose, R. (2015). Social media for social change: Social media political efficacy and activism in student activist groups. *Journal of Broadcasting & Electronic Media*, 59(3), 456–474. doi:10.1080/08838151.2015.1054998
- Verba, S., Lehman Schlozman, K., & Brady, H. E. (1995). *Voice and equality: Civic voluntarism in American politics*. Cambridge, MA: Harvard University Press.
- Watson, B. R., Peng, Z., & Lewis, S. C. (2019). Who will intervene to save news comments? Deviance and social control in communities of news commenters. *New Media & Society*, 21(8), 1840–1858. doi:10.1177/1461444819828328
- Wessler, H. (2008). Deliberativeness in political communication. In W. Donsbach (Ed.), *The international encyclopedia of communication* (pp. 1-6). doi:10.1002/9781405186407.wbiecd011
- Wilhelm, C., & Joeckel, S. (2019). Gendered morality and backlash effects in online discussions: An experimental study on how users respond to hate speech comments against women and sexual minorities. *Sex Roles*, 80(7–8), 381–392. doi:10.1007/s11199-018-0941-5
- Wilhelm, C., Ziegler, I., & Joeckel, S. (2019). Reporting hate comments: Investigating the effects of deviance characteristics, neutralization strategies and users' moral orientation. *Communication Research*. Advance online publication. doi:10.1177/0093650219855330
- Yeo, S. K., Su, L. Y. F., Scheufele, D. A., Brossard, D., Xenos, M. A., & Corley, E. A. (2017). The effect of comment moderation on perceived bias in science news. *Information, Communication & Society*, 22, 129–146. doi:10.1080/1369118X.2017.1356861