

Article

## Fueling Toxicity? Studying Deceitful Opinion Leaders and Behavioral Changes of Their Followers

Puck Guldemon<sup>1,\*</sup>, Andreu Casas Salleras<sup>2</sup>, and Mariken van der Velden<sup>2</sup>

<sup>1</sup> Strategic Communication Group, Wageningen University & Research, The Netherlands

<sup>2</sup> Department of Communication Science, Free University Amsterdam, The Netherlands

\* Corresponding author ([puck.guldemon@wur.nl](mailto:puck.guldemon@wur.nl))

Submitted: 3 May 2022 | Accepted: 26 August 2022 | Published: 30 December 2022

### Abstract

The spread of deceiving content on social media platforms is a growing concern amongst scholars, policymakers, and the public at large. We examine the extent to which influential users (i.e., “deceitful opinion leaders”) on Twitter engage in the spread of different types of deceiving content, thereby overcoming the compartmentalized state of the field. We introduce a theoretical concept and approach that puts these deceitful opinion leaders at the center, instead of the content they spread. Moreover, our study contributes to the understanding of the effects that these deceiving messages have on other Twitter users. For 5,574 users and 731,371 unique messages, we apply computational methods to study changes in messaging behavior after they started following a set of eight Dutch deceitful opinion leaders on Twitter during the Dutch 2021 election campaign. The results show that users apply more uncivil language, become more affectively polarized, and talk more about politics after following a deceitful opinion leader. Our results thereby underline that this small group of deceitful opinion leaders change the norms of conversation on these platforms. Hence, this accentuates the need for future research to study the literary concept of deceitful opinion leaders.

### Keywords

computational communication science; disinformation; opinion leaders; social media; the Netherlands; Twitter

### Issue

This article is part of the issue “Negative Politics: Leader Personality, Negative Campaigning, and the Oppositional Dynamics of Contemporary Politics” edited by Alessandro Nai (University of Amsterdam), Diego Garzia (University of Lausanne), Loes Aaldering (Free University Amsterdam), Frederico Ferreira da Silva (University of Lausanne), and Katjana Gattermann (University of Amsterdam).

© 2022 by the author(s); licensee Cogitatio (Lisbon, Portugal). This article is licensed under a Creative Commons Attribution 4.0 International License (CC BY).

### 1. Introduction

Social media (SM) platforms play a key role in our daily lives. People increasingly use SM to interact with friends and family, voice their opinions, consume news, and engage in politics (Popan et al., 2019; Spohr, 2017; Weeks et al., 2017). However, some information on SM is misleading, i.e., untrue, partly true, and potentially purposefully deceitful. This has been studied by scholars using the concepts of fake news (Egelhofer & Lecheler, 2019; Guess et al., 2019), disinformation (Bennett & Livingston, 2018; McKay & Tenove, 2021), rumors (DiFonzo & Bordia, 2007; Friggeri et al., 2014), or conspiracies (Douglas

et al., 2019; Sunstein & Vermeule, 2009), amongst other concepts. Academics have raised concerns, stating that deceitful content endangers democracy and society at large (Groshek & Koc-Michalska, 2017; Lazer et al., 2018; McKay & Tenove, 2021). For example, deceitful content has led to online discussions between SM users in which uncivil language is common, especially when these discussions are about polarizing political topics (Weeks & Gil de Zúñiga, 2021). This tone of voice, in turn, fuels toxicity on SM platforms (Kim et al., 2021). That is, uncivil language spills over to other SM users, which affects their attitudes towards those who are addressed in these messages, potentially leading to polarization.

Uncivil messages are those that contain curse words, are insulting, harassing, very dismissive towards others, racist or against a minority group, or are misogynistic, enabling a toxic sphere (Davidson et al., 2017; Theocharis et al., 2016). To remedy toxicity on their platforms, keeping them a healthy place for public debate, SM companies often remove users who spread deceitful content. This removal has fueled a societal debate about whether these actions are justified because such regulatory measures stand in contrast to the claims that SM would provide more equal opportunities for the free expression of political views than traditional media (Balkin, 2017). Hence, SM companies as private actors are engaged in regulating the “practical conditions of speech” in the digital space (Balkin, 2017). Yet, are those who spread deceitful content harmful to others? Currently, the empirical evidence on if and to what extent deceitful content harms other platform users is scarce. Therefore, we need a systematic study on disseminating a variety of types of deceitful content (e.g., fake news, conspiracies, rumors, and disinformation) and the effect thereof on other SM users.

We argue that the current state of the field aiming to understand the negative consequences of deceitful content is limited in three ways. First, previous research has been very compartmentalized. Scholars have studied different types of deceitful content in isolation (Weeks & Gil de Zúñiga, 2021). We, however, claim that when disseminating deceitful information, SM accounts spreading deceitful content often do not stick to just one type of deceitful information: They spread a variety of deceitful content throughout. Anecdotally, the now purged SM accounts of far-right radio show host Alex Jones show that he engages in conspiracies, as well as rumors and misleading information (Berr, 2019; Coaston, 2018; Haselton, 2019; Paul, 2019; Rosdorff, 2018). In our empirical analysis, we assess the validity of our claim that these kinds of salient accounts engage in the spread of different types of deceiving information. Thereby, this study meets and expands the work of Weeks and Gil de Zúñiga (2021), who call for research that goes beyond the mere distinction between different types of deceitful information. Furthermore, we build upon the work of Chadwick and Stanyer (2022), who theoretically argue for the need to have an overarching framework bridging the myriad of studies addressing deceitful content. We theorize and empirically demonstrate how various types of deceitful content are addressed, allowing us to gauge the harm of this content to other users and, thereby, to democracy and society at large. Second, existing research is focused on the type of content spread rather than on the SM accounts disseminating this information. If we aim to better understand the effect these salient accounts have on other users and, to some extent, whether the removal of accounts disseminating deceitful content is justified, we argue that not the content but the SM accounts should be at the center of analysis. We are interested in the negative effects that all these types of deceitful content

have on SM users, not just a particular type of deceitful content. Building upon the two-step flow of communication theory (Katz & Lazarsfeld, 1955), we introduce a theoretical concept that examines these salient accounts fueling SM platforms with toxicity by posting deceitful content, terming them “deceitful opinion leaders” (DOLs). Third, while there is plenty of existing knowledge about the overall prevalence and dissemination of deceitful content, we know little about the effects that DOLs have on their followers and other users on SM platforms. Scholars suggest that exposure to deceitful messages can have harmful consequences, such as adopting more uncivil behavior, and lead to increasing levels of affective polarization (Iyengar et al., 2012, 2019; Popan et al., 2019; Theocharis et al., 2016; Yarchi et al., 2021). We test whether exposure to deceitful information actually has these malicious effects.

To empirically assess the type of deceitful content DOLs spread as well as its effect on other SM users, we use an innovative research design that allows us to study the type of content DOLs disseminate and if this fuels toxicity on Twitter. Twitter is known as a key mainstream platform that allows us to collect the data needed to test our hypotheses. For a two-month period (March 2, 2021, till May 4, 2021), we tracked eight Dutch DOLs (Maurice de Hond, Lange Frans, Sietske Bergsma, Robert Jensen, Blackbox News, Wierd Duk, Cafe Weltschmerz, and Isa Kriens) and their followers. These DOLs are not an exhaustive nor representative list of DOLs in the Dutch Twittersphere. Yet, they are well known for engaging in the dissemination of deceitful information (e.g., see “YouTube verdedigt verwijderen account Lange Frans,” 2020), and thereby a most likely case to test our approach and theoretical concept. All DOLs have accounts with a high number of followers (i.e., more than 11,000), showing that these DOLs voice opinions that are valued and accepted by others. Moreover, DOLs often spread deceitful content about highly polarized and political issues. This results in an (online) public space fueled with toxicity and deceitful content (Bergmann, 2020). The collected messages of these DOLs allow us to assess the validity of our claim that DOLs engage in the spread of different types of deceiving content. For each day in the period under investigation, we monitored each DOL for if they had new followers ( $N = 32,245$ ). Subsequently, for each of these new followers, we collect the tweets they posted before and after they started following a DOL. Our analysis is two-fold. First, we look at the tweets posted by the DOLs and use content analysis to corroborate that they indeed engage in a wide variety of deceitful content, such as rumors and disinformation. Then, we look at the tweets sent by the new followers before and after and use computational methods to test the extent to which they become more politically engaged and post more uncivil and affectively polarized messages after following a DOL.

We show that, after starting to follow one of the eight DOLs in our sample, these users did increase their

number of political, uncivil, and affectively polarized tweets. The effects are statistically significant and of substantial magnitude. We observe stronger longer-term (30 days) than shorter-term (15 days) effects, although, after two weeks, their behavior starts reverting to levels similar to those before following the DOL. We also observe stronger effects for those who started following more than one DOL. Our results thereby underline that while there is a small group of DOLs, they do have a substantial effect on how other SM users behave on these SM platforms. To keep SM platforms a healthy forum for public debate, SM companies regulate what can be posted. Fueled by fear that the dissemination of deceitful information distorts a healthy public debate and, thereby, is detrimental to society, SM accounts engaging in this behavior are purged. Our results, however, demonstrate that following a DOL has a gateway effect: Not only are SM users adopting their norms of conversation (i.e., using more uncivil language), but they also introduce their SM followers to a view of politics that these followers feel more comfortable to engage in. This sheds important light on the question of how to regulate SM platforms so that they can maintain fostering public debate without endangering the democratic process of deliberation.

## 2. Deceitful Opinion Leaders on Social Media and Their Effects

Over the last decades, the media environment has changed drastically into a high-choice media environment (Van Aelst et al., 2017). This has affected the communication flow from the media to the masses. Many people receive news via SM through one of their online connections (Weeks et al., 2017). Hence, these connections function as a mediator between the media and the mass public. This process was first explained by Katz and Lazarsfeld (1955) as the two-step flow of communication theory, which acknowledges this process of person-to-person influence and calls these mediators opinion leaders. Those are people that are held in high esteem and whose opinions are valued and accepted by others (Bergström & Jervelycke Belfrage, 2018; Choi, 2015; Katz & Lazarsfeld, 1955). In the early days of mass media, opinion leaders received information from the media and shared that information with their network via (offline) personal interactions. In the digital age, this process is similar but takes place in an online environment: SM users seek out certain individual SM accounts for guidance and information (Choi, 2015). The information that SM users are exposed to depends on the opinions, interests, and behavior of their online connections (Bergström & Jervelycke Belfrage, 2018). These opinion leaders inform and thereby potentially shape the attitudes of less active recipients (Bergström & Jervelycke Belfrage, 2018; Carlson, 2019). Yet, they do not necessarily need to be message carriers for the greater good. In recent years, we have witnessed opinion leaders that

deliberately spread information that is untrue or deceiving, such as Alex Jones or Lange Frans in the Dutch context. Influential accounts that engage in this behavior we coin as DOLs. DOLs are defined as SM users (a) with a large number of followers and (b) who engage in the production and dissemination of at least one type of deceitful content to their audiences.

Why do people follow DOLs, and what is the effect thereof? Previous research demonstrates that most people are not necessarily engaged with politics, but do enjoy following entertaining content. As a by-product of seeking entertainment, politically inattentive individuals are exposed to information about political and societal issues (Baum, 2002). Social networks like Twitter provide increasing opportunities for people to be exposed to political content, even when using Twitter for different purposes, such as entertainment (Kim et al., 2013). DOLs typically post highly entertaining and engaging content, such as sarcastic or cynical comments. Hence, people are, in part, likely to follow them for entertainment value. A side-effect of following DOLs is that their followers are *incidentally exposed* (Bergström & Jervelycke Belfrage, 2018; Kim et al., 2013; Weeks et al., 2017) to political content—i.e., when DOLs tweet about societal, controversial, and political issues, their followers (and the followers of their followers via sharing patterns) see this content. The same dynamic holds for exposure to misleading information (Lazer et al., 2018; Stroud, 2008). We argue that DOLs have a key role in the information others receive, resulting in a high influence on what DOL followers talk about (Zaller, 1992). That is, the topics of conversation—i.e., the deceitful information about societal and political topics—likely spillover to the DOL followers, leading to the following hypothesis:

H1: After following DOLs, users will tweet more about politics than they did before following them.

Next to *what* DOLs talk about, *how* they speak about political topics is also likely to be carried over to their followers. According to Weeks and Gil de Zúñiga (2021), online political interactions are often uncivil. The highly emotional nature of SM platforms provides a “perfect storm” for the spread of deceiving and misleading content (Weeks & Gil de Zúñiga, 2021). DOLs often use inflammatory and uncivil rhetoric when discussing political topics or when referring to politicians (for an example, see Table 3). Due to anonymity, the threshold for uncivil behavior is lowered on SM platforms (Groshek & Koc-Michalska, 2017; Theocharis et al., 2016). In an SM environment, people tend to say and do things that they would not necessarily do when being in the offline world (Suler, 2004). Therefore, SM platforms facilitate this uncivil behavior online (Groshek & Koc-Michalska, 2017). This, in turn, results in the usage of more uncivil language, posing threats, hard criticism, and showing anger and hatred on SM platforms online, creating an online sphere rife with uncivil behavior (Suler, 2004;

Theocharis et al., 2016). Impolite and uncivil discourse on SM platforms has a poisonous and polarizing effect. When people are exposed to incivility, they are more likely to use incivility in their comments and messages (Gervais, 2015; Theocharis et al., 2016). This implies that those following DOLs, who are expected to use uncivil and inflammatory language, are more likely to mimic their rhetorical style, leading to the following hypothesis:

H2: After following DOLs, users will utilize more uncivil language.

Uncivil behavior on SM platforms reduces openness towards outgroups, as uncivil discourse has poisonous and polarizing effects (Groshek & Koc-Michalska, 2017; Theocharis et al., 2016). As mentioned above, DOLs often talk about political topics or politicians in an uncivil manner. They use an “us versus them” rhetoric when referring to the political elite. By doing so, they create an in-group (DOLs and their followers) and an out-group (the political elite and their followers). Based on the social identity theory (Tajfel & Turner, 1979), scholars have theorized and demonstrated that belonging to an in-group with a strong social identity leads to the disliking and disfavoring of out-groups (Harteveld, 2021; Iyengar et al., 2012, 2019). Online, this results in SM users following more like-minded accounts that fit within their in-group. This implies that once SM users follow a DOL, they are likely to be immersed in an online community of like-minded people, forming online homogenous networks (Barberá, 2015; Barberá et al., 2015; Shu et al., 2017). These homogeneous social networks reduce the tolerance for alternative worldviews and amplify affective polarization, resulting in division and animosity between different parties, individuals, or groups that hold opposite views on (political) topics (Iyengar et al., 2012, 2019; Lazer et al., 2018; Yarchi et al., 2021). Assuming that SM users following and mimicking a DOL can be seen as a united front (i.e., in-group), they are likely to view others as an out-group whom they

oppose. Thereby, they presumably contribute to rising hostility toward other societal groups. By following DOLs, we expect users to become less tolerant, hence more polarized, towards outgroups with different opinions and ideas. Therefore, we expect the following:

H3: After following DOLs, users will become more affectively polarized.

### 3. Data and Methods

We collected the following data to assess the extent to which DOLs engage in the dissemination of different kinds of deceitful content, as well as to test our three hypotheses about the effects that they have on the behavior of their followers. First, we selected a set of DOLs to study. Then, we tracked their SM behavior to explore the types of deceitful content they posted. In addition, we needed to track the SM behavior of their followers. Ideally, for clear identification, we wanted to track and study their behavior before versus after they started following a given DOL.

For a two-month period during the 2021 Dutch elections (March 2 through May 4, 2021), we studied the Twitter behavior of a convenience sample of eight well-known Dutch DOLs (for a detailed list, see Table 1) and those ordinary users who started following them during the period of analysis. Although these DOLs are not a representative nor comprehensive sample of all DOLs, they are among the most visible ones in the Dutch Twittersphere, and they are very suitable to conduct a proof-of-concept analysis to validate the theoretical concept and expectations put forward in this article. Future research should address the conditions under which the findings presented here extend to a larger and more comprehensive sample of DOLs. Despite this limitation, we believe the approach and analysis presented here contribute to building a better understanding of the actions of these types of opinion leaders and how they shape conversations on SM platforms.

**Table 1.** List of the eight DOLs we study.

Name	Twitter handle	Number of followers March 2, 2021	Number of followers May 5, 2021	Number of new followers analyzed (H1–H3)
Maurice de Hond	@mauricedehond	118,237	127,404 (+7.7%)	2,558
Wierd Duk	@wierdduk	84,617	90,403 (+6.8%)	1,101
Lange Frans	@langefrans	70,744	72,021 (+1.8%)	235
Robert Jensen	@robertjensen	52,686	56,462 (+7.1%)	588
Sietske Bergsma	@sbergsma	31,224	35,165 (+12.6%)	321
Café Weltschmerz	@cafeweltschmerz	17,062	17,754 (+4%)	32
Blck Bx	@blckbxnews	16,141	22,474 (+39.2%)	382
Isa Kriens	@isakriens	11,658	12,931 (+10.9%)	632
Total	—	—	32,245 (13,337 unique)	5,574

On the first day (March 2, 2021), we pulled the list of followers for each of these DOLs. We only include followers that have sent at least one tweet before to enable a comparison before and after these Twitter users started to follow a DOL. Then, every day (until May 4, 2021), we pulled the following additional information: the messages sent by the DOLs that day, the list of users who started following a given DOL that day, (up to) the last 3,200 messages sent by these new followers (to gather information about their posting behavior before following the particular DOL), and the messages posted that day by the new followers detected in previous days (to gather information about their posting behavior after they started following a particular DOL).

We use the collected data for two main purposes. First, we manually code the messages posted by the DOLs themselves for whether they contain fake news, disinformation, conspiracy, and/or rumors (non-mutually exclusive categories). The goal is to assess our claim that these DOLs engage in the dissemination of different types of deceitful content. As shown in Table 2, we rely on existing and validated definitions when coding for these four types of deceitful messages (see Part B in the Supplementary File for the codebook). For each DOL, 10 tweets were coded by two authors, resulting in 80 annotated tweets in total, leading to intercoder reliability values using Krippendorff's alpha of 0.99 for fake news, 0.98 for disinformation, 0.97 for conspiracy, and 0.93 for rumors.

Then, to test potential behavioral changes, we count the number of political (H1), uncivil (H2), and affectively polarized (H3) tweets that new followers posted during the days before versus the days after they started following the first DOL in our sample. We use two time windows for this before/after analysis, 15 and 30 days, to assess the robustness of the findings to this subjective cut-off. We collected data from 13,377 unique new followers for the DOLs in our sample. For clear identification, when testing our hypotheses, we will restrict our sample to (a) users who started following one of the DOLs after March 2, 2021 (for the previous followers, we do not know exactly the date they started following the DOL), (b) users for which we have collected their messages for the entire before and after time windows, and (c) users

who did not stop following the followed DOL during data collection (a total of 3,451 users started following one of the eight DOLs under analysis, but stopped following them before the end of data collection). Our final analytical sample includes a total of 5,574 followers (see Table 1) who sent a total of 731,371 tweets during the 30 days prior/after combined.

To count the number of political, uncivil, and affectively polarized tweets, we trained three machine-learning classifiers. First, we annotated 2,896, 5,242, and 855 for whether they were uncivil, political, and affectively polarized, respectively (binary categories). Table 4 provides an overview of the annotated messages per classifier. Messages were coded as uncivil if they were insulting, harassing, very dismissive towards others, racist or against a minority group, misogynist, or when they contain curse words (Davidson et al., 2017; Theocharis et al., 2016). Messages were coded as political if (a) a political party or organization was mentioned and/or (b) if messages touched on relevant policy issues. Finally, messages were coded as being affectively polarized if users showed dislike towards an opposing group (by naming them, tagging them, or mentioning them), such as a politician, political party, or societal group (e.g., conservatives/liberals, immigrants; see Part C in the Supplementary File for the codebook). One hundred Tweets were coded by two authors, leading to intercoder reliability values using Krippendorff's alpha of 0.86 for political tweets, 0.85 for uncivil language, and 0.87 for affectively polarizing language.

Since uncivil, political, and affectively polarized tweets are rare, to have as many true positives in our annotated set as possible, we used random sampling as well as active learning when selecting the cases to be annotated (Miller et al., 2020). Hence, the number of true positives in our annotated dataset is not really a reflection of the prevalence of these quantities in the overall dataset. Table 3 shows examples of the types of messages coded as political, uncivil, and affectively polarized.

Then we used the full corpus of annotated data to fine-tune three times the same transformer model (the Dutch version of BERT [de Vries et al., 2019]—bert-base-dutch-cased), one for each of the three (political, uncivil, and affectively polarized tweets) classifiers.

**Table 2.** Definitions of deceitful content used for coding.

Fake news	Has a journalistic format but is low in facticity (Egelhofer & Lecheler, 2019)
Disinformation	False information that is purposely spread to deceive people, seeking to amplify social divisions and distrust (Bennett & Livingston, 2018; McKay & Tenove, 2021)
Conspiracy	Efforts to explain events, practices, or secret plots that consist of two or more powerful actors acting in secret for their benefit and working towards a malevolent or unlawful goal against the common good (Douglas et al., 2019; Sunstein & Vermeule, 2009)
Rumor	Circulating information whose veracity status is yet to be verified at the time of spreading (DiFonzo & Bordia, 2007; Friggeri et al., 2014)

**Table 3.** Examples of political, uncivil, and affectively polarized tweets (translated from Dutch).

Political tweet	According to the left-wing opposition parties, the deal does not go far enough, while the PVV believes that the cabinet has caved in.
Uncivil tweet	@DDStandard she's ugly. she's stupid...she's not adding anything. just a hopeless nigger who also tries to shout something...and nobody listens. she will never become someone like Pim Fortuyn...Sylvana cannot even stand in his shadow
Affective polarized tweet	RT I didn't think much of the left wing voters, but voting for fucking Sigrid Al Qaq-Kaag is like selling your soul to Europe...

Deep transformer models such as BERT have been shown to improve machine text classification in many domains, including political and communication science (Terechshenko et al., 2020). In each case, we used 20% of the annotated data to create a completely untouched validation set. Then, we split 70/30 of the remaining data into a train and test set. We used the train set to estimate model fit and update the model weights at each training iteration and the test split to assess out-of-sample performance and to decide when to stop training the model further. We stopped the training when the test loss did not improve for three complete iterations. We trained each model three times, using a different train/test split each time (three-fold cross-validation). Finally, we assessed out-of-sample accuracy on the untouched validation set (which remained constant across the three folds).

In Table 4, we report the performance of each model based on this three-fold cross-validation conducted on the validation set. The uncivil and political classifiers perform very well: Overall accuracy, as well as precision and recall, are very high; and precision and recall are very similar, indicating that in the rare cases in which a classifier makes the wrong prediction, it is equally likely to misclassify messages that are (vs. are not) uncivil/political. The performance of the affective polarization classifier is slightly lower—high accuracy (83%) but slightly lower levels of precision (65%) and recall (71%)—but the classifier is highly balanced (similar levels of precision and recall). We have no reason to believe that there is any systematic error for any of the classifiers. So, any remaining noise would mean that we are conducting conservative tests of our hypotheses.

Finally, we use these classifiers to predict whether the rest of the unlabeled messages posted by the new 5,574 DOL followers are political, uncivil, and affectively polarized, and to count the number of political/uncivil/polarizing tweets sent the 30 days before and

the 30 days after starting to follow the first DOL in our sample.

#### 4. Results

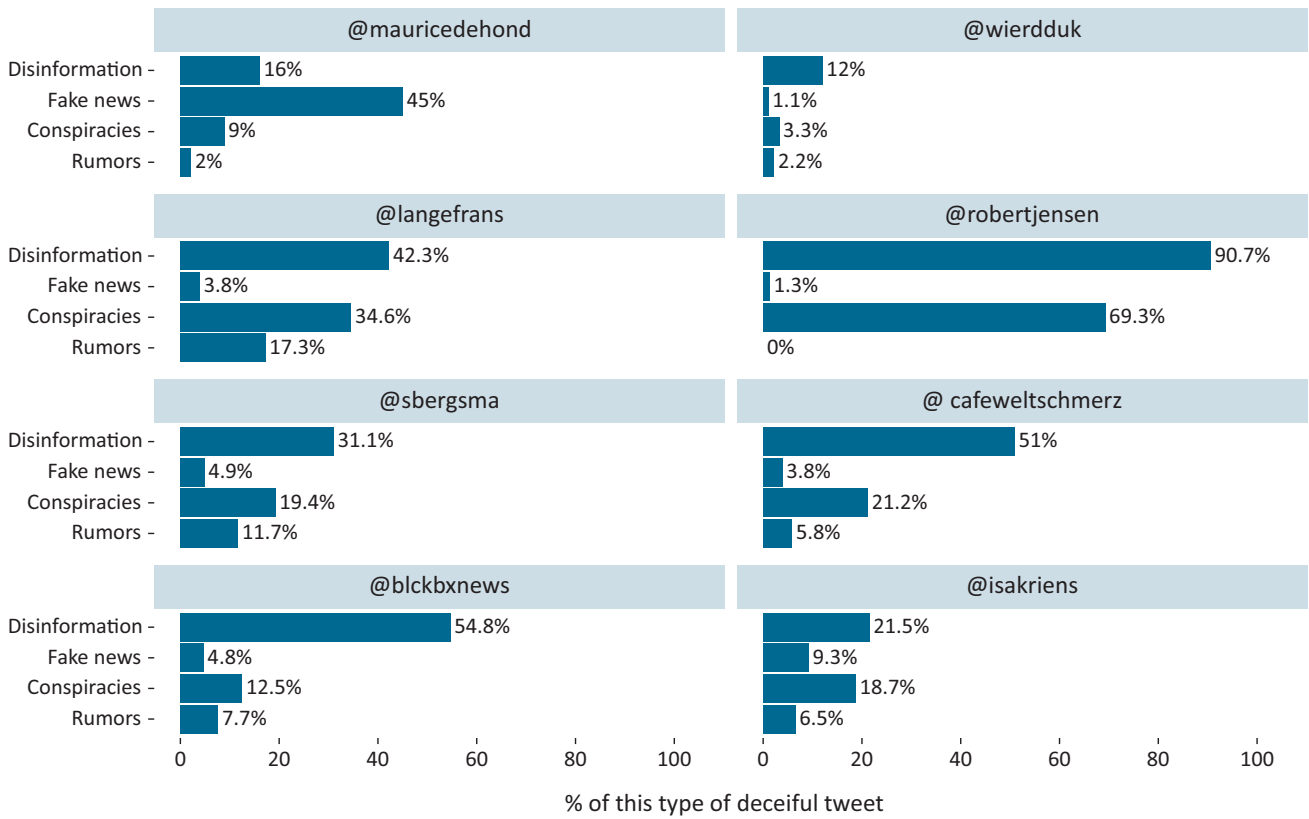
We begin by assessing whether DOLs indeed engaged in the dissemination of many types of deceitful content (e.g., fake news, disinformation, conspiracies, and rumors). We then move to test our hypotheses regarding the behavior of new followers.

In Figure 1, we study the distribution of the deceitful content that was spread by each DOL during the period of analysis. In line with our theoretical framework, the figure illustrates that all DOLs engage (to some extent) in the dissemination of all types of deceitful content under scrutiny, from fake news to conspiracies and rumors. For example, except for Robert Jensen, the remaining DOLs posted at least one message containing each of the deceitful typologies under study. Although sometimes they have a clearly preferred deceitful category (e.g., 45% of Maurice de Hond's tweets spread fake news, and 54.8% of Blck Bx's messages promoted disinformation), they also engage in the spread of other kinds of deceitful content quite often (e.g., 16% and 9% of Maurice de Hond's tweets contained disinformation and conspiracies, respectively; and 12.5% and 7.7% of Blck Bx's messages had conspiracies and rumors in them). These results align with our argument that the main goal of these actors is to inject toxicity into online environments and that each type of deceitful content is simply one of many tools in the toolbelt of DOLs. In addition, the results emphasize that a user-centric (rather than, or in combination with, a content-centric and compartmentalized) analysis is needed to have a clearer understanding of the spread of deceitful content on SM and its effects.

To test H1, H2, and H3, we turn to the set of new followers for which we had collected enough information

**Table 4.** Three-fold cross-validated performance of three BERT classifiers predicting binary outcomes: Political, uncivil, and affectively polarized tweets.

Classifier	N annotated	True positives	Accuracy	Precision	Recall
Political	5,242	59%	86%	94%	86%
Uncivil	2,896	39%	86%	83%	80%
Affectively Polarized	855	35%	83%	65%	71%

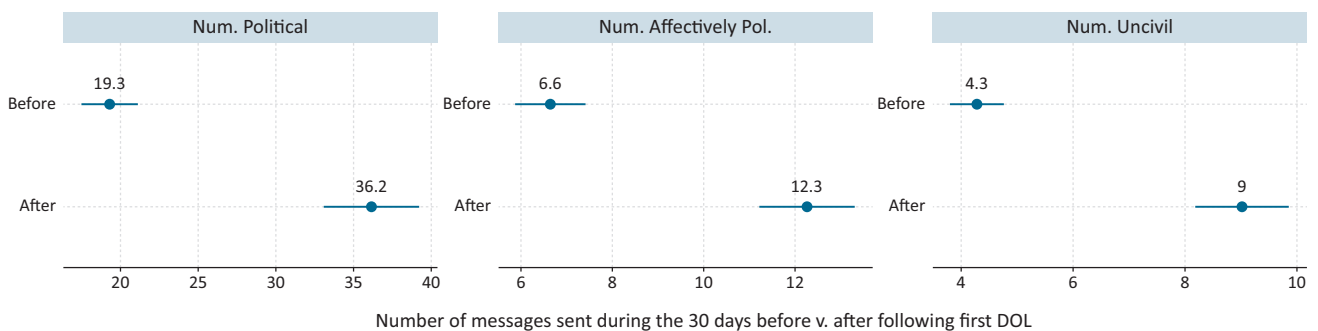


**Figure 1.** The percentage of tweets sent by the DOLs under analysis that contain different types of deceitful content.

to explore a potential change in behavior after following the first DOL in our sample (N = 5,574). In Figure 2, we show the average number of political, uncivil, and affectively polarized tweets these users sent during the 30 days before (vs. after) following the first DOL. We see stark differences across the board. The users were more politically engaged (sending 36.2 political tweets in the 30 days after vs. 19.3 political tweets in the 30 days prior), more uncivil (9 vs. 4.3 uncivil tweets), and affectively polarized (12.3 vs. 6.6 polarizing tweets).

Given that we collected the data during an election period, we wanted to control for whether a user started following a DOL before the election day (as users may have been more likely to discuss politics during the *after* time window). Hence, we created the vari-

able Campaign Post Days, which accounts for the number of *post* 15/30 days that overlapped with the electoral campaign (so the number of days between the day a user started following the first DOL and election day, March 17). This variable is 0 for those who started following a deceitful opinion leader after March 17. As specified in Model 1, for a clearer test of our hypotheses we use linear models predicting the difference ( $Y_{post} - Y_{pre}$ ) for three outcomes of interest (number of uncivil, affectively polarized, and political tweets) as a function of the mentioned control variable Campaign Post Days. For each of these linear models, the intercept parameter ( $\alpha$ ) provides information about the average difference in messaging behavior between the post and pre-difference after accounting for the control variable.



**Figure 2.** Average number of political, uncivil, and affectively polarized tweets (plus 95% confidence interval) sent during the 30 days before (vs. after) following the first DOL.

Model 1, the model specification used to test H1, H2, and H3, is as follows:

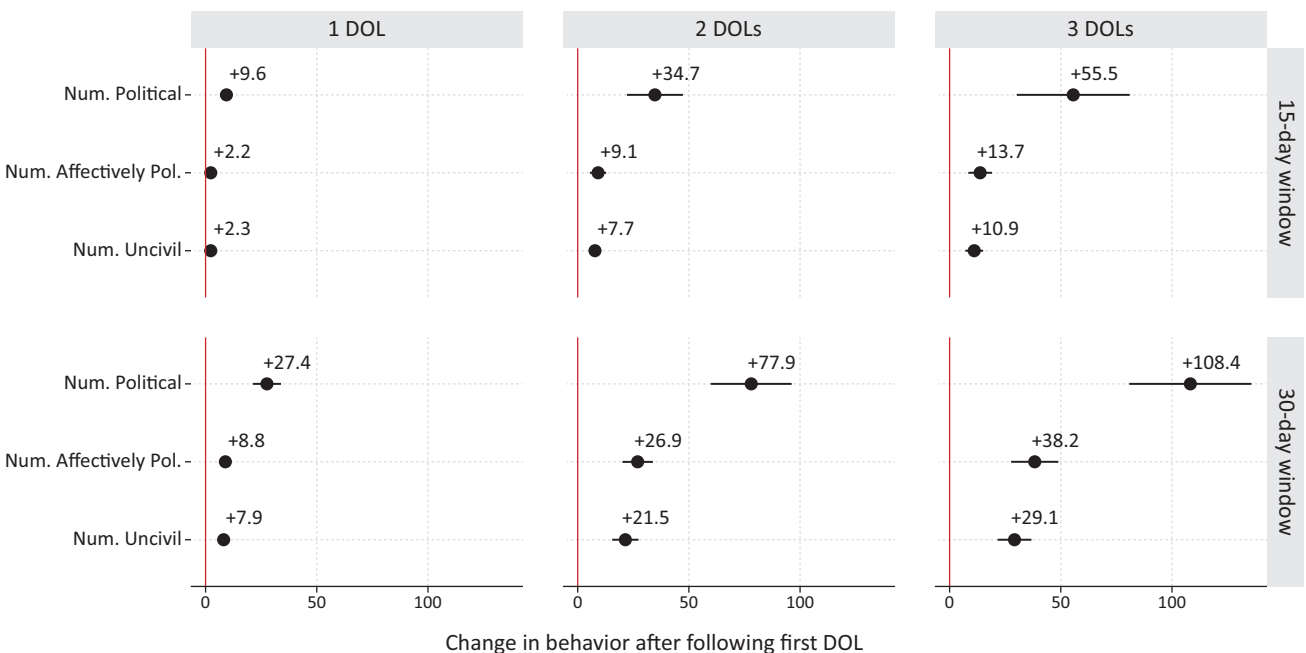
$$(Y_{post} - Y_{pre}) = \alpha + \beta_{CampaignPosDays} + \epsilon$$

In Figure 3, we report the  $\alpha$  coefficient for several linear models. For each of the four outcomes of interest, we ran six models with the same specification (i.e., Model 1), where we varied the time window to calculate the post/pre periods (15- and 30-day windows) and the number of DOLs the user followed within the 15/30 days after following the first opinion leader. In the first column (1 DOL), we include all the users in our sample ( $N = 5,574/3,891$ ), and in the other columns we estimate the models using only those users who followed at least a second DOL (2 DOLs) within the next 15/30 days ( $N = 1,336/1,014$ ), and at least a third DOL (3 DOLs;  $N = 555/421$ ); i.e., each analysis includes the number of unique Twitter users that meet the criteria. These variations allow us to disentangle differential effects across time (whether we observe stronger effects when comparing 15 vs. 30 days), and across different levels of engagement (e.g., users who decided to follow more than one of the DOLs in our sample).

We find strong support for our three hypotheses. Across the board, we see an increase in the number of political, uncivil, and affectively polarized tweets. All estimates presented in Figure 3 are statistically significant at the conventional 0.05 level. We observe the mildest effects among those who only followed one of the DOLs in our sample. But even among those, we observe a substantial change in behavior, particularly when we compare the behavior during the 30 days after (vs. before) following the DOL. On average, these users sent 27.4 more political messages, 8.8 more affectively polarized messages, and 7.9 more uncivil messages.

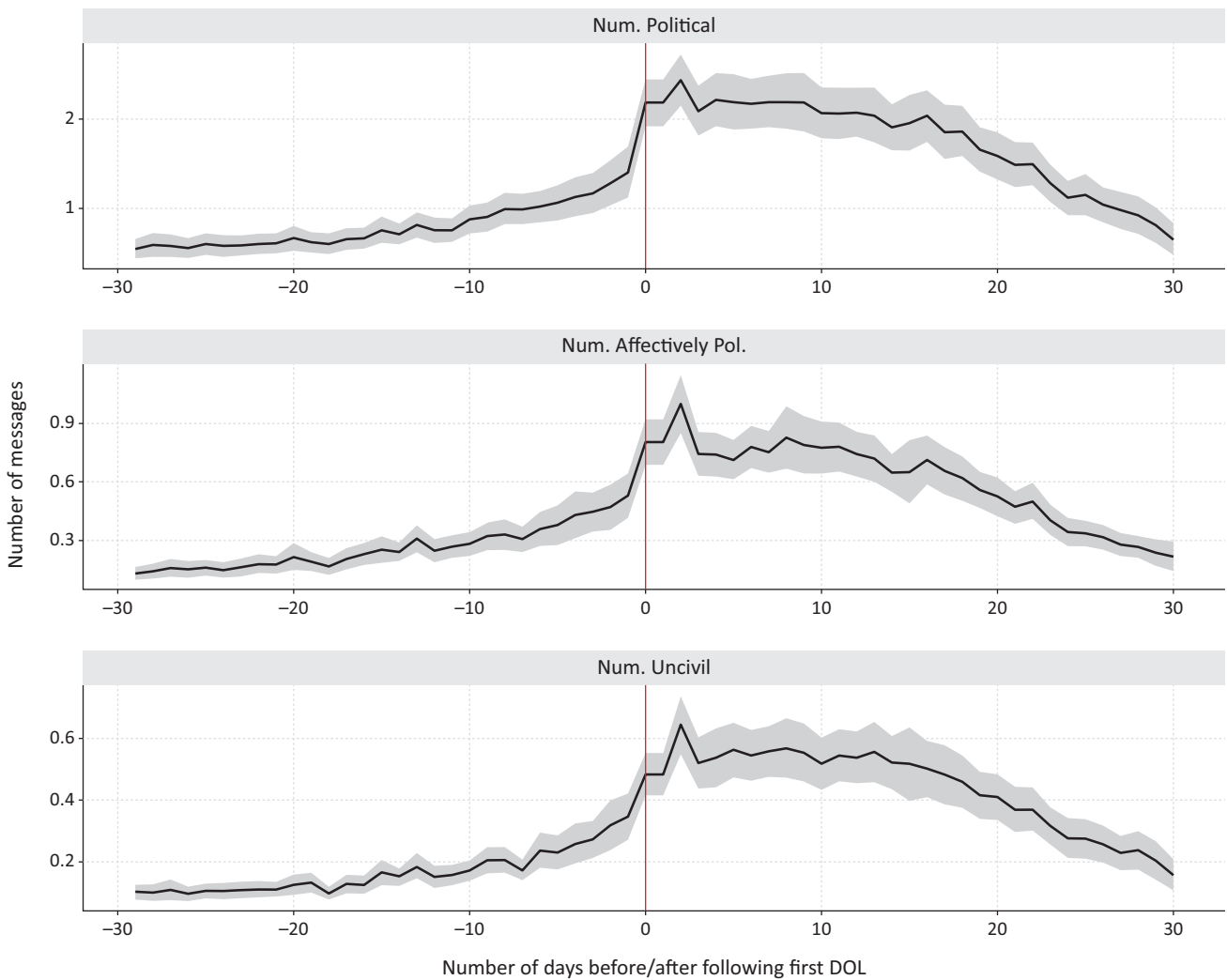
tively polarized messages, and 7.9 more uncivil ones. We observe the strongest effects among those who followed a second and a third DOL during the 30 days after following the first DOL in our sample ( $N = 421$ ). On average, they radically sent more political (+108.4), affectively polarized (+38.2), and uncivil (+29.1) messages. These findings are not driven by the new followers of one particular DOL, but reflect a general pattern observed across the followers of the different DOLs in our sample (see Appendix D in the Supplementary File). In addition, this behavior change cannot be simply explained by these users retweeting messages originally posted by the DOLs they started following (see Appendix A in the Supplementary File). On average (95% confidence intervals included), only 0.8% (0.6–1%) of the political tweets, 0.5% (0.3–0.7%) of the uncivil messages, and 0.5% (0.3–0.7%) of the affectively polarized tweets they sent during the 30 days after following the first DOLs are retweets of that DOL.

For a more detailed picture, in Figure 4, we explore the functional form of these effects. The figure shows the average number of political, uncivil, and affectively polarized tweets (+95% confidence intervals) the users in our sample sent each of the 30 days before and after following the first DOL. Figure 4 is standardized, so the exact date of day 0 differs across users, depending on when they started following the DOL. We observed a slight upper trend right before they started following the DOL. This indicates that at least some users already started shaping their behavior before day 0. This could be because they may have already been exposed to some tweets from these DOLs via retweets from their networks, or some factor motivated them to change their behavior and potentially seek these kinds of opinion



**Figure 3.** Coefficients (+95% confidence intervals) from linear models estimating a change in behavior after following one, two, and three DOLs.





**Figure 4.** Average number of political, uncivil, and affectively polarized tweets sent by followers of DOLs, during the 30 days before and after following the first DOL (N = 3,891).

leaders. Then, we observe a clear jump at the moment the users started following the first DOL. The number of overall tweets and the uncivil, affectively polarized, and political ones remained high for about 15 to 20 days. After that period, the behavior of the users gradually reverted to their levels of activity before following the DOL. The patterns described in Figure 4 clearly point to these DOLs playing a crucial role in the radicalization of online environments. Independently of what motivated these users to start following these DOLs—whether it was a very intentional decision or because of incidental exposure via retweets from one’s networks—we observe stark and substantive changes in behavior that contribute to increasing levels of toxicity and incivility on the SM platform.

### 5. Conclusion and Discussion

This article tackles three shortcomings of existing literature studying the dissemination of deceitful content. First, existing literature is very compartmentalized,

as it mostly focuses on one type of deceitful content (e.g., conspiracies, fake news, or misleading information). We show that salient SM accounts engage in the spread of all sorts of deceitful content throughout. Each type of deceitful content that they disseminate is just one of many tools in their toolbox. Second, we lack an overarching approach that puts these influential SM users at the center instead of the content that they spread. We do so by putting forward a new theoretical concept: “deceitful opinion leaders.” Third, this study contributes to the understanding of the individual-level effects that these types of deceiving messages have on other SM users. We show that after following a DOL on Twitter, significant behavioral changes start to occur amongst their followers: users send more political, uncivil, and affectively polarized messages. For example, on average, the analyzed users sent around 28 political tweets, 8 uncivil tweets, and 9 affectively polarized tweets more during the 30 days after following a DOL, compared to the 30 days prior. These behavioral changes seem to gradually revert to their levels of activity before following the

DOL. Although at the individual level these behavioral changes do not last long, at the aggregate level these effects have a substantive impact: DOLs gather new followers every day, meaning that these behavioral effects are constantly occurring, having a longer-lasting effect on the behavior and norms of conversation on Twitter.

Although this article adds important results to existing literature, it is not without limitations. This article provides a first aim in studying the effects of DOLs on SM platforms. There are other influential DOLs who were not included in this research. Moreover, all the DOLs in this study are Dutch. Hence, this study only focuses on the Dutch SM landscape. Furthermore, this study only considers Twitter, while DOLs are active on many platforms. To assess the generalizability of the effects that DOLs have on other SM users, future studies should aim to address additional factors that influence these findings, such as platform affordances and the level of radicalization of a platform. We expect the work presented here to inspire future work focusing on a more comprehensive and representative sample of DOLs from different contexts on different platforms, to provide further insights into the conditions under which these opinion leaders shape our online environments. Despite these limitations, this research finds valid and important results that show significant individual-level effects from following DOLs who engage in the spread of deceitful content online. Even though research finds that only a small proportion of SM users spread deceitful content per se (e.g., Guess et al., 2019), the spread of deceitful content via SM leads to substantial effects on other users on the platform.

The results of this study provide a first look into the distribution of the spread of deceitful content by DOLs and the individual-level effects that DOLs have on their followers. Importantly, this study adds to the empirical evidence of the effects of deceitful content on SM users. The findings of this study add to existing literary knowledge of the consequences of deceitful content in online environments. In addition, the results of this study provide empirical evidence to the societal debate on whether these influential SM users should be removed to maintain a healthy forum for public debate. Removing DOLs from Twitter would reduce toxicity on the platform. However, doing so might have negative effects if DOLs move on to other platforms to spread their deceiving content. This might result in higher levels of radicalization and polarization. Especially on Telegram, which is known to have a high number of users that support conspiracy theories. Furthermore, these findings underline that a small group of DOLs change the norms of conversation on SM platforms. Hence, this accentuates the need for future research to study the literary concept of DOLs.

### Acknowledgments

This research is funded by an NWA Small Project Grant (“When Words Do the Trick: The Effect of Rhetorical

Strategy on the Perceived Legitimacy of Political Decision-Making”) awarded to Dr. Mariken van der Velden and Dr. Andreu Casas Salleras. We thank the members of the Political Communication Group at the Free University Amsterdam, and the participants of the workshop “Negative Politics: Leader Personality, Negative Campaigning, and the Oppositional Dynamics of Contemporary Politics” and the panel “How the Type of Media Use Affects Belief in Conspiracy Theories and Misinformation” at ICA’22 for their thoughtful comments and suggestions on early versions of this manuscript.

### Conflict of Interests

The authors declare no conflict of interests.

### Supplementary Material

Supplementary material for this article is available online in the format provided by the author (unedited).

### References

- Balkin, J. M. (2017). Free speech in the algorithmic society: Big data, private governance, and new school speech regulation. *UCDL Review*, 51, 1149–1210.
- Barberá, P. (2015). Birds of the same feather tweet together: Bayesian ideal point estimation using Twitter data. *Political Analysis*, 23(1), 76–91. <https://doi.org/10.1093/pan/mpu011>
- Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., & Bonneau, R. (2015). Tweeting from left to right: Is online political communication more than an echo chamber? *Psychological Science*, 26(10), 1531–1542. <https://doi.org/10.1177/0956797615594620>
- Baum, M. A. (2002). Sex, lies, and war: How soft news brings foreign policy to the inattentive public. *American Political Science Review*, 96(1), 91–109. <https://doi.org/10.1017/S0003055402004252>
- Bennett, W. L., & Livingston, S. (2018). The disinformation order: Disruptive communication and the decline of democratic institutions. *European Journal of Communication*, 33(2), 122–139. <https://doi.org/10.1177/0267323118760317>
- Bergmann, E. (2020). Populism and the politics of misinformation. *Safundi*, 21(3), 251–265. <https://doi.org/10.1080/17533171.2020.1783086>
- Bergström, A., & Jervelycke Belfrage, M. (2018). News in social media: Incidental consumption and the role of opinion leaders. *Digital Journalism*, 6(5), 583–598. <https://doi.org/10.1080/21670811.2018.1423625>
- Berr, J. (2019, June 25). Conspiracy theorist Alex Jones evades crackdown on social media sites. *Forbes*. <https://www.forbes.com/sites/jonathanberr/2019/06/24/conspiracy-theorist-alex-jones-evades-crackdown-on-social-media-sites/?sh=1809755a4f87>
- Carlson, T. N. (2019). Through the grapevine: Infor-

- mational consequences of interpersonal political communication. *American Political Science Review*, 113(2), 325–339. <https://doi.org/10.1017/S000305541900008X>
- Chadwick, A., & Stanyer, J. (2022). Deception as a bridging concept in the study of disinformation, misinformation, and misperceptions: Toward a holistic framework. *Communication Theory*, 32(1), 1–24. <https://doi.org/10.1093/ct/qtab019>
- Choi, S. (2015). The two-step flow of communication in Twitter-based public forums. *Social Science Computer Review*, 33(6), 696–711. <https://doi.org/10.1177/0894439314556599>
- Coaston, J. (2018, August 6). Alex Jones banned from YouTube, Facebook, and Apple, explained. *Vox*. <https://www.vox.com/2018/8/6/17655658/alex-jones-facebook-youtube-conspiracy-theories>
- Davidson, T., Macy, M., Warmesley, D., & Weber, I. (2017). *Automated hate speech detection and the problem of offensive language*. arXiv. <https://arxiv.org/abs/1703.04009>
- de Vries, W., van Cranenburgh, A., Bisazza, A., Caselli, T., van Noord, G., & Nissim, M. (2019). *BERTje: A Dutch BERT model*. arXiv. <https://doi.org/10.48550/ARXIV.1912.09582>
- DiFonzo, N., & Bordia, P. (2007). Rumor, gossip and urban legends. *Diogenes*, 54(1), 19–35. <https://doi.org/10.1177/0392192107073433>
- Douglas, K. M., Uscinski, J. E., Sutton, R. M., Cichocka, A., Nefes, T., Ang, C. S., & Deravi, F. (2019). Understanding conspiracy theories. *Political Psychology*, 40(S1), 3–35. <https://doi.org/10.1111/pops.12568>
- Egelhofer, J. L., & Lecheler, S. (2019). Fake news as a two-dimensional phenomenon: A framework and research agenda. *Annals of the International Communication Association*, 43(2), 97–116. <https://doi.org/10.1080/23808985.2019.1602782>
- Friggeri, A., Adamic, L., Eckles, D., & Cheng, J. (2014). Rumor cascades. In E. Adar & P. Resnick (Eds.), *Proceedings of the International AAAI Conference on Web and Social Media* (pp. 101–110). Association for the Advancement of Artificial Intelligence.
- Gervais, B. T. (2015). Incivility online: Affective and behavioral reactions to uncivil political posts in a web-based experiment. *Journal of Information Technology & Politics*, 12(2), 167–185. <https://doi.org/10.1080/19331681.2014.997416>
- Groshek, J., & Koc-Michalska, K. (2017). Helping populism win? Social media use, filter bubbles, and support for populist presidential candidates in the 2016 US election campaign. *Information, Communication & Society*, 20(9), 1389–1407. <https://doi.org/10.1080/1369118X.2017.1329334>
- Guess, A., Nagler, J., & Tucker, J. (2019). Less than you think: Prevalence and predictors of fake news dissemination on Facebook. *Science Advances*, 5(1), Article eaau4586. <https://doi.org/10.1126/sciadv.aau4586>
- Harteveld, E. (2021). Fragmented foes: Affective polarization in the multiparty context of the Netherlands. *Electoral Studies*, 71, Article 102332. <https://doi.org/10.1016/j.electstud.2021.102332>
- Haselton, T. (2019, May 3). Alex Jones was banned from Facebook, but an hour later he was back on Facebook livestreaming. *CNBC*. <https://www.cnn.com/2019/05/02/alex-jones-banned-from-facebook-but-hes-already-back.html>
- Iyengar, S., Lelkes, Y., Levendusky, M., Malhotra, N., & Westwood, S. J. (2019). The origins and consequences of affective polarization in the United States. *Annual Review of Political Science*, 22(1), 129–146. <https://doi.org/10.1146/annurev-polisci-051117-073034>
- Iyengar, S., Sood, G., & Lelkes, Y. (2012). Affect, not ideology. *Public Opinion Quarterly*, 76(3), 405–431. <https://doi.org/10.1093/poq/nfs038>
- Katz, E., & Lazarsfeld, P. F. (1955). *Personal influence: The part played by people in the flow of mass communications*. Free Press.
- Kim, J. W., Guess, A., Nyhan, B., & Reifler, J. (2021). The distorting prism of social media: How self-selection and exposure to incivility fuel online comment toxicity. *Journal of Communication*, 71(6), 922–946. <https://doi.org/10.1093/joc/jqab034>
- Kim, Y., Chen, H.-T., & Gil de Zúñiga, H. (2013). Stumbling upon news on the Internet: Effects of incidental news exposure and relative entertainment use on political engagement. *Computers in Human Behavior*, 29(6), 2607–2614. <https://doi.org/10.1016/j.chb.2013.06.005>
- Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., Metzger, M. J., Nyhan, B., Pennycook, G., Rothschild, D., Schudson, M., Sloman, S. A., Sunstein, C. R., Thorson, E. A., Watts, D. J., & Zittrain, J. L. (2018). The science of fake news. *Science*, 359(6380), 1094–1096. <https://doi.org/10.1126/science.aao2998>
- McKay, S., & Tenove, C. (2021). Disinformation as a threat to deliberative democracy. *Political Research Quarterly*, 74(3), 703–717. <https://doi.org/10.1177/1065912920938143>
- Miller, B., Linder, F., & Mebane, W. R. (2020). Active learning approaches for labeling text: Review and assessment of the performance of active learning approaches. *Political Analysis*, 28(4), 532–551. <https://doi.org/10.1017/pan.2020.4>
- Paul, K. (2019, May 3). Facebook bans Alex Jones, other extremist figures. *Reuters*. <https://www.reuters.com/article/us-facebook-extremists-usa-idUSKCN1S82D7>
- Popan, J. R., Coursey, L., Acosta, J., & Kenworthy, J. (2019). Testing the effects of incivility during internet political discussion on perceptions of rational argument and evaluations of a political outgroup. *Computers in Human Behavior*, 96, 123–132. <https://doi.org/10.1016/j.chb.2019.02.017>
- Rosdorff, M. (2018, August 8). Complotdenker Alex

- Jones verbannen van bijna alle sociale media [Conspiracy theorist Alex Jones banned from almost all social media]. *EenVandaag*. <https://eenvandaag.avrotros.nl/item/complotdenker-alex-jones-verbannen-van-bijna-alle-sociale-media>
- Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter*, 19(1), 22–36. <https://doi.org/10.1145/3137597.3137600>
- Spohr, D. (2017). Fake news and ideological polarization: Filter bubbles and selective exposure on social media. *Business Information Review*, 34(3), 150–160. <https://doi.org/10.1177/0266382117722446>
- Stroud, N. J. (2008). Media use and political predispositions: Revisiting the concept of selective exposure. *Political Behavior*, 30(3), 341–366. <https://doi.org/10.1007/s11109-007-9050-9>
- Suler, J. (2004). The online disinhibition effect. *CyberPsychology & Behavior*, 7(3), 321–326. <https://doi.org/10.1089/1094931041291295>
- Sunstein, C. R., & Vermeule, A. (2009). Conspiracy theories: Causes and cures. *Journal of Political Philosophy*, 17(2), 202–227. <https://doi.org/10.1111/j.1467-9760.2008.00325.x>
- Tajfel, H., & Turner, J. C. (1979). An integrative theory of intergroup conflict. In M. J. Hatch & M. Schultz (Eds.), *Organizational identity: A reader* (pp. 55–65). Oxford University Press.
- Terechshenko, Z., Linder, F., Padmakumar, V., Liu, F., Nagler, J., Tucker, J. A., & Bonneau, R. (2020). *A comparison of methods in political science text classification: Transfer learning language models for politics*. SSRN. <http://doi.org/10.2139/ssrn.3724644>
- Theocharis, Y., Barberá, P., Fazekas, Z., Popa, S. A., & Parrot, O. (2016). A bad workman blames his tweets: The consequences of citizens' uncivil Twitter use when interacting with party candidates. *Journal of Communication*, 66(6), 1007–1031. <https://doi.org/10.1111/jcom.12259>
- Van Aelst, P., Strömbäck, J., Aalberg, T., Esser, F., de Vreese, C., Matthes, J., Hopmann, D., Salgado, S., Hubé, N., Stępińska, A., Papathanassopoulos, S., Berganza, R., Legnante, G., Reinemann, C., Sheaffer, T., & 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. <https://doi.org/10.1080/23808985.2017.1288551>
- Weeks, B. E., & Gil de Zúñiga, H. (2021). What's next? Six observations for the future of political misinformation research. *American Behavioral Scientist*, 65(2), 277–289. <https://doi.org/10.1177/0002764219878236>
- Weeks, B. E., Lane, D. S., Kim, D. H., Lee, S. S., & Kwak, N. (2017). Incidental exposure, selective exposure, and political information sharing: Integrating online exposure patterns and expression on social media. *Journal of Computer-Mediated Communication*, 22(6), 363–379. <https://doi.org/10.1111/jcc4.12199>
- Yarchi, M., Baden, C., & Kligler-Vilenchik, N. (2021). Political polarization on the digital sphere: A cross-platform, over-time analysis of interactional, positional, and affective polarization on social media. *Political Communication*, 38(1/2), 98–139. <https://doi.org/10.1080/10584609.2020.1785067>
- YouTube verdedigt verwijderen account Lange Frans: "Richtlijnen meermaals geschonden" [YouTube defends the removal of Lange Frans account: "Guidelines have been violated several times"]. (2020, October 21). AD. <https://www.ad.nl/show/youtube-verdedigt-verwijderen-account-lange-frans-richtlijnen-meermaals-geschonden~a5606d7d>
- Zaller, J. (1992). *The nature and origins of mass opinion*. Cambridge University Press.

## About the Authors



**Puck Guldmond** is a PhD candidate at the Strategic Communication Group at Wageningen University & Research. Her research interests include data-driven campaigning, political influence and behavior on social media, computational social sciences, and populist rhetoric in political communication.



**Andreu Casas Salleras** is an assistant professor at the Department of Communication Science at the Vrije Universiteit Amsterdam and a faculty associate at the Center for Social Media and Politics at New York University. His research interests encompass the areas of political communication, public policy processes, and computational social sciences. He is particularly interested in how social movements and interest groups influence the political agenda and the decision-making process in the current media environment.



**Mariken van der Velden** is an associate professor of political communication at the Vrije Universiteit Amsterdam. She was a postdoctoral researcher at the Institute of Political Science at the University of Zurich and visiting researcher at the Departments of Political Science at the University of North Carolina, Chapel Hill, and the University of Oxford. Her research interest comprises the areas of political communication, political behavior, and computational social science.