

# Computational observation

*Challenges and opportunities of automated observation within algorithmically curated media environments using a browser plug-in*

Mario Haim & Angela Nienierza

CCR 1 (1): 79–102

DOI: 10.5117/CCR2019.1.004.HAIM

## Abstract

A lot of modern media use is guided by algorithmic curation, a phenomenon that is in desperate need of empirical observation, but for which adequate methodological tools are largely missing. To fill this gap, computational observation offers a novel approach—the unobtrusive and automated collection of information encountered within algorithmically curated media environments by means of a browser plug-in. In contrast to prior methodological approaches, browser plug-ins allow for reliable capture and repetitive analysis of both content and context at the point of the actual user encounter. After discussing the technological, ethical, and practical considerations relevant to this automated solution, we present our open-source browser plug-in as an element in an adequate multi-method design, along with potential links to panel surveys and content analysis. Ultimately, we present a proof-of-concept study in the realm of news exposure on Facebook; we successfully deployed the plug-in to Chrome and Firefox, and we combined the computational observation with a two-wave panel survey. Although this study suffered from severe recruitment difficulties, the results indicate that the methodological setup is reliable and ready to implement for data collection within a variety of studies on media use and media effects.

**Keywords:** data collection, computational methods, social media, observation, browser plug-in, Facebook, media use

The prevalence of Internet use has significantly changed the way people encounter information. Leaving aside changing media habits and routines, algorithms, in particular, have become increasingly important for filtering the overwhelming amount of information with which users are confronted online. Such algorithms present filtered selections of information, tailored to what ought to be most beneficial to users or to usage turnout for the relevant providers.

Yet, the individualized nature of information encounters is not reflected in much of today's media research. For example, content analyses may collect news articles directly from an outlet's website, thus implicitly assuming that all users encounter the same news content. Such studies therefore ignore the uniqueness of an individual user's content selection as well as the individual contexts of that content. While these shortcomings might be due to research-economic reasons, they nevertheless interfere with presumed micro-level media effects within algorithmically curated media environments (e.g., van Aelst et al., 2017; Zuiderveen Borgesius et al., 2016).

To study the increasingly important encounters with algorithmically curated information, researchers need adequate tools: that is, reliable and valid measures of the extent to which users encounter information and the context of the content to which they are exposed (e.g., Freelon, 2018). Unfortunately, not much is known about what content is presented in which context within algorithmically curated media environments, since intermediaries, such as Facebook or Google, act behind closed doors and hardly ever cooperate with independent research projects. Computational modes of enquiry seem to be a promising way to fill this methodological gap, but so far they have not been applied or evaluated to any great extent.

This paper therefore sets out to review the possibilities of computational observation, the unobtrusive and automated collection of information encountered in the course of a user's online behavior by means of a browser plug-in. In contrast to other methodological approaches, browser plug-ins allow for reliable capture and repetitive analysis of both content and context at the moment of the actual user encounter. Computational observation is thus a valid and reliable approach to the investigation of algorithmically curated media environments.

The objective of this paper is threefold. First, we describe the methodological challenges in assessing user encounters with information in algorithmically curated environments and compare the methodological approaches currently available. Second, we present a novel approach to unobtrusive observation of how people actually use such environments, in the form of a browser plug-in; we describe the challenges and benefits

of this method and discuss the central technological, ethical, and practical considerations. Third, we present a proof-of-concept study of news use within the algorithmically curated media environment of Facebook, combining a two-wave quantitative survey with a computational observation of participants' news use by means of the browser plug-in. The paper concludes with a discussion of the advantages and disadvantages of computational observation.

## 1. Prior Methodological Approaches

In the investigation of user encounters with information in algorithmically curated environments, three streams of methodological approaches seem applicable. First, studies might concentrate on the users by conducting surveys. Second, studies might focus on the algorithmically curated environment by collecting the available materials and subjecting them to content analysis. Third, studies might use observation techniques to investigate the point of contact between users and content.

## 2. Surveys

The traditional approach to measuring media use has been to ask respondents post hoc about their use. However, this kind of measurement has long been suspected of generating systematic error, as people tend to overestimate their exposure times (Haenschen, 2019; Price & Zaller, 1993; Prior, 2009; Scharkow, 2016). While previous research has shown that, for example, open-ended questions reduce such overreporting (Guess, 2015), post-hoc surveys provide no reliable way to capture individually encountered context. That is, although respondents might remember using certain news outlets, post-hoc surveys usually obscure information on various contexts, such as who shared the respective news headline or how many “likes” it had.

To address this shortcoming, individuals have also been surveyed using experience sampling methods. In this research procedure, respondents are asked to provide self-reports on occasions selected at random for close-up investigation of recently encountered media content (e.g., Reinecke & Hofmann, 2016; Struckmann & Karnowski, 2016). For example, Struckmann and Karnowski (2016) used mobile text messages to contact participants and ask them to report on their news consumption during that particular

day. While such surveys “in situ” (Reinecke & Hofmann, 2016, p. 456) can be quite burdensome to participants, they are more valid for the context-dependent investigation of algorithmic content curation than the usual post-hoc surveys, as they do not rely as heavily on the users’ memory.

Despite this improvement, both survey types remain of limited applicability to algorithmically curated media environments, for two main reasons. First, they are limited to content that the respondents have perceived; that is, if algorithms filter and show content but respondents do not take notice of it, conclusions about algorithmic curation can only be speculative. Second, surveys are dependent on respondents classifying content in the way the researchers intended; that is, they depend for example, on the respondents’ understanding of “politics” matching the researchers’ own definition of the term (Guess, Munger, Nagler, & Tucker, 2018).

### 2.1. Content Analyses

Instead of asking individual users about their media environment, one can look at the variety of content available. The easiest and most traditional way is to analyze the variety of available originating sources (e.g., news sites, politicians’ profile pages). Yet, while such a procedure has obvious advantages (e.g., it is not time-critical), it also falls short of providing any indication about algorithmic curation. Instead, it assumes that all users encounter the same content and ignores individual contexts.

Alternatively, a sample of algorithmically provided content could be collected, for example by adopting the exploratory approach of modeling different types of users and their online behavior (e.g., Haim, Graefe, & Brosius, 2016). On this approach, human profiles could be simulated for regular collection of algorithmically curated content; this increases the reliability of the data by allowing the researcher to control for various influences that are difficult to hold constant in real life. However, the main disadvantage of this method is that it produces an artificial environment, which limits the external validity of the results.

### 2.2. Observations

To achieve greater external validity, observational studies seem appropriate. Combined with self-confrontation interviews, they allow for insights into both media use and individually encountered context (e.g., Kümpel, 2018). However, because of the highly labor-intensive nature of this method, the number of participants to be observed and the collection of data are limited. Furthermore, there is a risk that participants might be influenced by the observation.

Less obtrusive observational approaches include log-file analyses (e.g., Scharkow, 2016), tracking manually uploaded data provided by various platforms (e.g., Menchen-Trevino, 2016), and passive tracking through third-party panels (e.g., Revilla, Ochoa, & Loewe, 2017). Yet, while these approaches allow for large-scale data collection, they do not provide insights into individually encountered contexts. That is, although they typically include the URLs of content that has been depicted and/or visited, they do not capture popularity cues, for example, or commentary that was presented along with the actual content.

### 3. Computational Observation Using a Browser Plug-In

To overcome these pitfalls, a promising approach to the observation of users within algorithmically curated media environments is by means of a browser plug-in. A browser plug-in allows for the unobtrusive observation of content, contexts, and actual user behavior, regardless of whether the content was encountered by users or simply shown without their notice (e.g., loaded but not scrolled over). Moreover, it allows the collection of user interactions to identify whether the content that was shown was eventually acted upon. A browser plug-in is also capable of collecting observational data on a large scale, as it depicts a standardized and thus quantitative way of data collection, thereby allowing measures to be collected from an appropriately high number of users. Consequently, quantitative statistical analyses are possible across different users, browsers, and settings, and even across different algorithmically curated media environments.

Despite these advantages, a browser plug-in is not yet a well-established approach to online observation, but rather an explorative exception within the growing field of computational communication research (for notable exceptions, see Menchen-Trevino, 2016; Möller, van de Velde, Merten, & Puschmann, 2019). As such, it involves a plethora of thus far undiscussed methodological challenges, including technological and ethical considerations, as well as assumptions about the actual research application.

In accordance with these considerations, we have developed and made available an open-source browser plug-in for computational observation.<sup>1</sup> After installation and user registration, the plug-in repeatedly contacts a centralized server to retrieve a configuration ruleset regarding which information to collect. Once retrieved, the plug-in unobtrusively collects information depicted within the browser. By collecting actually depicted information, the plug-in is almost independent of the platform under

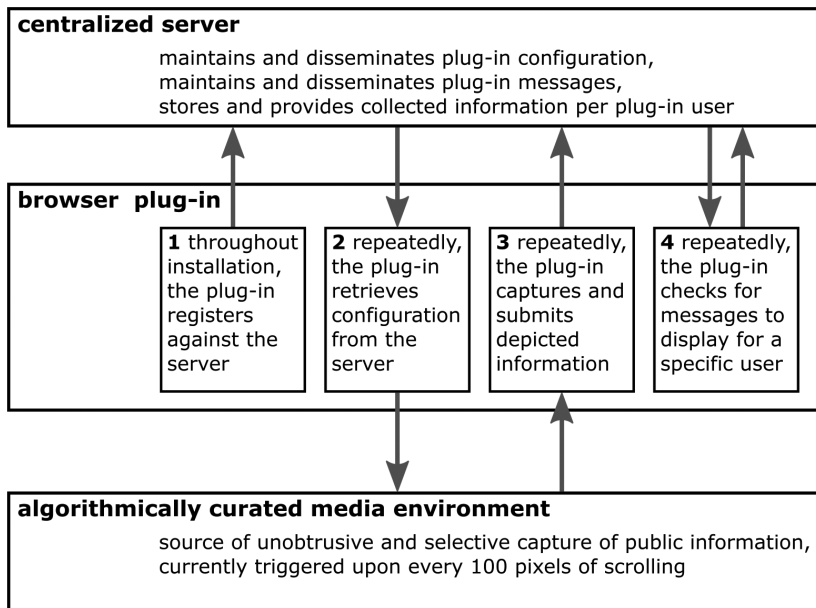


Figure 1. Basic functionality of the presented browser plug-in

investigation (e.g., encryption); instead, it only requires a specification of which information to collect. Upon data collection, the plug-in then submits this information to the server. The plug-in is also capable of showing messages from the server to the user, for example to invite users to an accompanying survey (Figure 1).

### 3.1. Technological Considerations

Technologically, two main elements are necessary for a browser plug-in to observe and collect content, context, and online user behavior unobtrusively. First, the browser plug-in itself needs to be developed separately for each targeted browser. Second, a centralized server needs to be set up to process data collected through installed plug-ins. This server can be addressed by plug-ins independently of the targeted browsers. In addition to collecting and processing observed data, the centralized server may be used to control the installed plug-ins; that is, it can orchestrate how the plug-ins collect content and which content they collect, as well as send messages to users at any time.

To combine plug-ins and server, a communication standard needs to be defined. Depending on the extent to which the server is able to orchestrate

plug-ins, this communication standard must, as a minimum, specify the format for submitting the content collected; it should also include messages to the users and a configuration ruleset on how to collect particular types of content. Importantly, because of the observed content's potential to obstruct privacy regulations, a communication standard should also adhere to a certain level of security. That is, the communication standard should (1) prevent unauthorized clients from accessing information, (2) ensure the credibility of submitted content, and (3) secure a certain level of anonymity. Our open-source plug-in therefore relies on hash-based message authentication (HMAC with SHA1 encryption).

**Browser plug-in.** The browser plug-in needs to be developed separately for different browsers. All the more common browsers (e.g., Apple Safari, Google Chrome, Microsoft Edge, Microsoft Internet Explorer, Mozilla Firefox, Opera) allow third-party plug-ins to be installed within their desktop versions.

For the development of browser plug-ins, some browsers follow their own plug-in technologies, platforms, and guidelines. These browsers thus require unique plug-in development, whereas other browsers have agreed on a quasi standard for core functionality in plug-in development. Plug-ins following the so-called “WebExtensions” (or “Extension API”) standard allow multi-browser deployment at almost no additional expense. This quasi standard has been supported by Google Chrome, Mozilla Firefox (from version 57), and Opera, and it has recently been announced as easily adaptable by Microsoft Edge (Microsoft, 2017). While this plug-in support is thus independent of the operating system being used, Apple Safari and Microsoft Internet Explorer have not adhered to WebExtensions. Browser features available to WebExtensions are limited to basic functionality supported by every browser, such as interaction with online content, display of messages, and invisible communication with third-party servers. In contrast, browser-specific features, such as sorting window tabs or handling address-bar interactions, are unavailable to plug-ins based on WebExtensions.

We developed our browser plug-in following the WebExtensions standard, since the major functionality is covered by this quasi standard. Our browser plug-in can (1) observe a browser's depiction of online content, (2) identify individual information within the layout of a website under observation, (3) capture user interactions such as scrolling, clicking, or typing, and (4) submit the information collected to the centralized server. To complement this computational observation with surveys, the browser plug-in can also (5) show messages with hypertext content (e.g., forms, links).

Upon installation, the plug-in asks the user for an anonymized identification code and a password as login credentials so that participants can later inspect and potentially delete the information that has been collected about them. Once installed, the plug-in repeatedly pulls mappings of the information to be observed and their corresponding selectors from the centralized server. This procedure allows for centralized maintenance of the observation data within the potentially changing layout of an algorithmically curated media environment. Computational observation is then based on CSS3 selectors that are applied to the information environment on each instance of scrolling activity to capture not only initially loaded but also subsequently reloaded information. Finally, the plug-in can be served with individual messages by the centralized server. In relation to a participant's anonymized identification code, these messages can be distributed individually and personalized for each user.

**Centralized server.** In contrast to the browser plug-in, the centralized server does not need browser-specific development. Rather, it can easily follow the constraints of a REST API (Representational State Transfer Application Programming Interface), a uniform interface for machine-to-machine communication that handles client-server calls individually and thus neither requires constant communication channels nor relies on sophisticated client management.

For the development of the centralized server in the form of a REST API, various libraries are available for a plethora of programming languages. Importantly, for more web-friendly programming languages, such as Java, PHP, Python, or even C++, lots of open-source frameworks for easy setup of REST APIs are maintained, including Spring (Java), Laravel (PHP), Flask (Python), and Wt (C++). These languages and frameworks are also capable of communicating with a variety of database systems, the ultimate selection of which is not of the utmost importance to the API's functionality. While each language and framework comes with its own caveat, no prerequisites apply other than the API's need to be in line with the development of the plug-in. To this end, we have included JSON schemes in the plug-in source repositories.<sup>2</sup>

### 3.2. Ethical Considerations

From an ethical perspective, computational observation potentially deals with personal information while providing multiple opportunities for data infringement in the course of data processing. Consequently, such data requires profound protection, and users must be provided with transparency



and full control over their data, as is stipulated, for instance, within the European Union's General Data Protection Regulation (GDPR). In following these regulations within the process of computational observation, three points in time are critical: before, during, and after observation of online user behavior.

First, before any user behavior is observed, users must give their informed consent to being observed: "consent should be given by a clear affirmative act establishing a freely given, specific, informed and unambiguous indication of the data subject's agreement to the processing of personal data relating to him or her" (European Parliament & Council of the European Union, 2016, sec. 32). Hence, at the time of installation, a browser plug-in should inform the user about the specific information that will be collected and should request explicit consent to the intended observation. It is therefore important to consider that only the installing user gives their consent; if a plug-in is installed on a shared or public computer, we must ensure that the online behavior of the consenting user only will be observed.

Second, during data collection, the GDPR (European Parliament & Council of the European Union, 2016, sec. 25) asks for "data protection by design and by default" and refers to "appropriate technical and organizational measures." Both plug-in and server should incorporate anonymization wherever possible, transmit data using encryption and ensuring data integrity, and ideally avoid collecting information that is capable of identifying individuals. Usernames or e-mail addresses may be pseudonymized using hashing algorithms, client-server communication can make use of HTTP over TLS and include various forms of integrity-checking authorization, such as bearer tokens or any other hash-based message authentication.

Third, after data collection, the user must be able to gain knowledge of all the data collected about him or her (European Parliament & Council of the European Union, 2016, sec. 12), and a user has "the right to withdraw his or her consent at any time" (European Parliament & Council of the European Union, 2016, sec. 7, par. 3). Consequently, users should be able to pause computational observation manually. Moreover, data should be accessible to users, as a minimum upon request, and data structures should not depend on individually collected information. That is, users possess a "right to be forgotten" (European Parliament & Council of the European Union, 2016, sec. 17), which might lead to deletion of observed information but should not affect information collected for other users. While our plug-in allows observation to be paused manually and information to be pseudonymized at the point of data collection, there are additional ethical

considerations that focus on data provision and deletion. As such, a centralized server should take care of these aspects.

Finally, depending on the algorithmically curated media environment under investigation, automatic collection of large amounts of data may be prohibited by the providers' terms of use (e.g., Freelon, 2018). Yet, while this paper cannot serve as a legal advisor in any sense, some scholars consider the research context as exceptional, at least within the European Union; this makes the collection of *public* information a gray area subject to case-by-case evaluation (e.g., Bodó et al., 2018; Diakopoulos, 2014).

### 3.3. Practical Research Considerations

Observation-based research requires further practical considerations to be taken into account, three of which are particularly noteworthy.

First, the dissemination of the plug-in. While browser plug-ins can be installed locally through so-called "developer modes," a more practical and widespread method is installation through official plug-in stores. Requirements for uploading one's own plug-ins into stores vary; for example, Google charges developers a small fee to become eligible for plug-in dissemination, and an automated code review is conducted. Likewise, Mozilla asks for plug-in code to be open source, otherwise a manual code review is initiated, which might take up to two weeks. In addition, and despite the shared quasi standard of plug-in development, every plug-in store asks for individual meta information and promotional materials (e.g., screenshots, logos, descriptions, product websites). On installation, plug-ins then require users to permit certain privileges when requested, such as accessing depicted information for a specified list of domains and communicating with a third-party server. Consequently, disseminated plug-ins should request privileges for the algorithmically curated media environment(s) under investigation while our provided open-source plug-in is capable of universal application.

Second, although browser plug-ins allow observation of actual depictions of online content and user interactions with that content, they are susceptible to layout changes on the websites that are being observed. For example, if the number of Facebook "likes" is to be captured, CSS or XPath selectors are an appropriate technological specification for identifying the number of "likes." Such selectors represent a hierarchical path along an HTML layout. If a website under observation changes its HTML layout, these selectors might miss the mark until they have been updated. Consequently, if such selectors are hard-coded into the browser plug-in, layout changes require updating the entire plug-in,

which might take some time owing to plug-in store restrictions. An alternative is for plug-ins to be served dynamically by the centralized server with mappings of the information to be retrieved and the corresponding selectors. If there is an HTML layout change to the website under observation, the plug-ins do not need to adapt but only to request updated mappings, which is a quicker process. Our plug-in does this by relying on a JSON-encapsulated set of rules to apply CSS3 selectors to the information depicted.

Third, data protection regulations and general privacy concerns may have fostered skepticism toward data-collecting projects, hence dampening participant recruitment. Possible solutions include providing a transparent message about the research purpose and providing incentives to participate (e.g., money, personal insights into online behavior).

#### 4. Proof-of-Concept Study

To put the computational observation approach to the test, we conducted a proof-of-concept study. While this study does not meet academic criteria for sample sizes or representativeness, it nevertheless demonstrates the applicability of our methodological approach. It also demonstrates how to employ the plug-in and make effective use of it within (future) research endeavors.

We conducted our proof-of-concept study in the realm of online news consumption within social networking sites, as such sites have become increasingly important for accessing news from around the world. For example, according to the 2018 Reuters Digital News Report (Newman, Fletcher, Kalogeropoulos, Levy, & Nielsen, 2018), Facebook is the most widely used social media platform for accessing news, and more than one third (36%) of digital news consumers reported using Facebook for news. However, in studying this increasingly important source of news, researchers face three major methodological challenges: (1) fragmentation of media use, (2) fragmentation of media sources, and (3) personalization.

(1) Respondents' media repertoires in today's high-choice media environment vary widely because of structural, sociodemographic, and habitual factors (e.g., Kim, 2014; Taneja, Webster, Malthouse, & Ksiazek, 2012; Wolfsfeld, Yarchi, & Samuel-Azran, 2016). In addition, repertoire fragmentation may be fueled by the large variety of content available on social networking sites, where respondents might have different views on what counts as "journalistic news." This fragmentation of media use limits the

reliability of the available survey data and hence impedes research into media use (Hasebrink & Popp, 2006; Webster & Ksiazek, 2012).

(2) Given that social networking sites act as intermediaries that point to existing third-party content, their users encounter news from a variety of originating sources. These diffusion mechanisms have led news publishers, as well as alternative and partisan sites, to adjust their distribution practices toward strategies that are optimized for social media (Newman et al., 2018; Newman, Fletcher, Levy, & Nielsen, 2016). Along with this variety of *originating* sources, the selection of encountered news also varies with regard to the content's *disseminating* sources. Social networking sites allow users to discover news that has previously been shared by the originating news source itself, by befriended people, or by following other routes, such as appearances on subscribed pages. Users may therefore also encounter news incidentally that they would not otherwise actively seek (Ahmadi & Wohn, 2018). From a methodological perspective, this fragmentation of media sources blends and mixes media brands for users and thus further increases the risk of incorrect self-reporting of media use (e.g., Kalogeropoulos & Newman, 2017; Newman et al., 2016).

(3) The abundance of content available on social networking sites is hardly manageable for individual users. Therefore, personalization mechanisms filter information at a per-user level (Bozdag, 2013; Webster & Ksiazek, 2012). This has led to variation in the contexts within which users encounter news, ranging from different popularity cues (Haim, Kümpel, & Brosius, 2018) to different accompanying user comments (Anderson, Brossard, Scheufele, Xenos, & Ladwig, 2014). Although, surprisingly, many Facebook users seem to be unaware that their personal feed is algorithmically curated (Eslami et al., 2015), personalization highlights the need for empirical investigation of the contexts in which users encounter news.

Our proof-of-concept study aims to take account of these major challenges for research into media use within social networking sites. It combines the method of computational observation with a two-wave panel survey to collect and compare the news that was actually encountered with the news use that was reported post hoc. The environment under investigation is Facebook, arguably one of the world's most influential algorithmically curated media environments. This multi-method approach allows us to estimate the importance of Facebook within respondents' media repertoires, to look at the actual content respondents encountered in their news feeds, and to investigate the contexts in which the news content was depicted, allowing to account for fragmentation of media use, fragmentation of media sources, and personalization. Due to the study's demonstrational

character, we refrain from formulating hypotheses. Instead, this proof-of-concept study aims (1) to provide an estimate of respondents' willingness to participate, (2) to report on the characteristics of sessions, posts, and interactions, and (3) to compare observed media use with self-reported indicators of media use.

#### 4.1. Method

We developed a centralized server and deployed the browser plug-in with privileges requested for the Facebook domain under the name of "FBforschung.de Browser-Plugin" to the plug-in stores of Google<sup>3</sup> and Mozilla<sup>4</sup>. The plug-in, depicted as a sleeping gray owl, blended in neatly with other plug-in symbols on the browser's top menu bar. When a participant opened Facebook, the owl unobtrusively "woke up" and its color turned Facebook-blue. Clicking the owl symbol allowed participants to visit the project's website to examine and potentially delete any information that had been captured.

**Procedure.** Participants were asked to answer an initial survey at  $t$ . This was intended to lower barriers to participation, as users are expected to be more proficient in answering surveys than in installing browser plug-ins. The survey asked for the participant's media repertoire and various indications of Facebook use, both generally and with particular regard to news.

If participants responded that they used Facebook at least sometimes on a desktop computer inside either Google Chrome or Mozilla Firefox, the survey presented them with an easy installation manual and referred them to the relevant plug-in store. On installation, the plug-in (1) informed participants about the data to be collected and required their consent, (2) asked for an anonymized identification code consisting of the first letter of the user's mother's first name and other similarly unobtrusive characters, and (3) requested participants to log in to Facebook to tie an anonymized MD5 representation of the alphanumeric Facebook user identifier to the plug-in's anonymized identification code. This last step was designed to ensure that only the consenting participant's Facebook news feed was observed; any other users logging into Facebook within the browser in which the plug-in was installed were dismissed.

Two weeks after installation, the browser plug-in received an order from the centralized server to present participants with a message asking them to answer the second survey, which was linked and included the anonymized identification code. Since the appearance of this message was individually generated two weeks after installation, rather than scheduled to a specific

date, the time of this survey is referred to as  $t_2$ . All data collection took place in October and November of 2018.

**Observational data.** Every visit to Facebook's main news feed website (i.e., facebook.com after login) was captured and subsequently referred to as a "session." Meta information stored for each session included a time stamp, the version numbers of both the plug-in and the regularly pulled information-selector mappings, the browser (along with its version and operating system), the browser's language, and the number of pixels that the user scrolled past.

Within each session, the browser plug-in collected "posts," visually demarcated elements within the news feed. To respect the privacy and data protection of third-party users (e.g., Facebook friends of our study participants), the plug-in collected only posts that were public, as indicated by Facebook with a small globe symbol. The data captured for each post included the initial publication time stamp, textual content and imagery, embedded links, the names of originating and disseminating users, groups, or pages, the numbers of "likes" and other popularity cues, and (where available) the textual hint as to why a post was visible (e.g., "a friend shared this"). Moreover, the number of visible comments below the post, the post's position as counted ordinally from the top of the news feed, and the post's position counted from the top in pixels, were also stored.

Finally, the plug-in captured user interactions with a post. The interactions observed included clicking on a post's link and "liking" a post. While this case study is primarily intended to provide a proof of concept, the observation of interactions could easily be extended to, for example, commenting on a post, expressing other reactions (e.g., "love," "haha"), and sharing a post.

**Survey data.** For this proof-of-concept study, the primary aims were to verify the technological capability of the research design and to compare observed and self-reported news use within Facebook. In the survey at  $t_2$ , we therefore asked participants how much they used various social networking sites. We also asked them to estimate the average proportions in their news feed of (a) private posts, (b) news posts, (c) advertisements and commercial posts, (d) posts from subscribed celebrities or pages, and (e) other posts. In order to get a sense of the importance of Facebook within each participant's media repertoire, we also asked how much they used other news sources and how important those other sources were, following the questions in the Reuters Digital News Report (Newman et al., 2018).

In accordance with the German Longitudinal Election Study (Roßteutscher et al., 2018), participants were also asked about any political actions they had taken online during the previous year. As in the European Values Study (Gedeshi et al., 2010), we asked individuals to report their levels of trust in various institutions, such as the media, police, political parties, and Facebook. The survey ended with sociodemographic questions, the construction of the anonymized identification code, and an installation manual for the plug-in.

The second survey, at  $t_2$ , asked participants how much they had used Facebook in the previous two weeks (the period of observation). As at  $t_1$ , we asked them to estimate the proportions of different kinds of posts over the previous two weeks so that we could compare their reports with our observations. The survey at  $t_2$  also asked participants to tell us about the importance for them of Facebook as a source of political information in the previous two weeks, their encounters with news from a variety of outlets, and their perceived incidental news exposure to (a) news from sources that they do not usually use, (b) news from sources that they did not subscribe to, (c) news that contradicted their own opinions, and (d) news that would not usually be of interest to them. These questions were designed to indicate how representative the observation period was for each individual.

**Participants.** Participants were recruited in three waves by means of invitation letters highlighting the innovative character of the study and its potential to support academic endeavors in shedding more light on algorithmically curated media environments. First, we tried to find willing students on Facebook by posting our request into several local user groups; however, this did not have a major recruitment effect. Second, we were able to send exactly 500 invitation emails through the SoSci panel, a non-commercial online panel for convenience sampling of German-speaking respondents (Leiner, 2014). Respondents on the SoSci panel have opted to be invited to academic studies up to twice a year. They can thus be regarded as research-savvy and open to arguments highlighting the importance of academic endeavors against the interests of private companies. While this round of invitations yielded some respondents, willingness to participate (and in some cases suitability) remained low. Hence, and third, we sent out a further 300 invitation emails through the SoSci panel. As an incentive, each respondent was eligible to take part in a lottery for one of five vouchers, valued at 25 euros, for an online shop of the winner's choosing.

## 4.2. Results

**Participation.** Despite these recruitment measures, only 8 participants installed the plug-in and completed both questionnaires. Such a small sample does not allow for any inferential statistical analyses. In addition, 15 participants installed the plug-in but did not answer the second survey, resulting in a total of 23 “plug-in participants” for whom we were able to gain observational data. Importantly, the survey at  $t_1$  was started 166 times, but only 104 questionnaires were completed. According to the SoSci panel committee, this response rate of 13 per cent is lower than their usual rates of 20 to 25 per cent; Sax and colleagues (2003) reported that for web-based surveys, response rates between 17 and 20 per cent are to be expected. While we can only speculate on the reasons for the low participation rate in this case, the necessity of installing a plug-in, mentioned in the recruitment e-mail, might have discouraged people from participating. In addition to this initial lack of motivation, potential participants also dropped out later in the process, because, for example, they did not use Facebook at all ( $n = 18$ ), they only used it on mobile devices ( $n = 3$ ), or they used it within unsupported browsers ( $n = 26$ ). Finally, 28 participants who had responded to the survey at  $t_1$  refused to install the plug-in (for undisclosed reasons).

Among the 23 plug-in participants, the average age was  $M = 43.2$  years ( $SD = 17.9$ ). 14 participants were female. 12 had a university degree, and 9 participants had a higher education entrance qualification. 18 participants installed the plug-in on Firefox, 5 of whom also completed the survey at  $t_2$ ; another 5 used the Chrome plug-in, and 3 of them also completed the survey at  $t_2$ . For all 23 participants, a total of 6,809 Facebook sessions were collected, consisting of 43,410 posts and 691 observed interactions. Despite the prominent display of information on how to do so, no participant deleted any session, post, or interaction from their data collection. Participants spent an average of  $M = 5.4$  minutes ( $SD = 3.3$ ) completing the survey at  $t_1$  and  $M = 4.6$  minutes ( $SD = 1.2$ ) completing the survey at  $t_2$ .

**Sessions, posts, and interactions.** Within the two-week period of observation, participants visited Facebook  $M = 296$  times on average ( $SD = 411.3$ ). While this suggests rather intensive use of Facebook, each manual refresh of the Facebook website counts as a new session, technically speaking. While we do not know much about the actual usage behavior of users on Facebook in this regard, the high number of sessions suggests that at least some participants repeatedly refreshed the Facebook page, and this renders the measurement of Facebook visits somewhat ambiguous. It is



cumbersome to try to identify a page refresh. As an estimate, one could assume that if two sessions from one plug-in were created within, for example, five minutes, this might be considered a page refresh rather than a new Facebook session. For the available dataset, such a highly artificial five-minute assumption would decrease the average number of sessions (i.e., five-minute visiting windows) down to  $M = 54$  ( $SD = 68.2$ ). This would be in line with our finding that for 3,234 sessions (48%), users did not scroll at all. On this hypothetical assumption, in-depth investigation of individual participants suggests a pattern of refreshing the Facebook website that is fairly evenly distributed across almost all participants. Sticking to the original definition of a session is therefore not expected to introduce systematic bias.

Sessions were created through almost the whole course of a day (Figure 2). Peaks can clearly be identified in the evening, especially between five o'clock in the afternoon and midnight: a total of 3,862 (57%) sessions were created in that time range. This may be due to the plug-in being limited to desktop computers and excluding mobile devices.

Within these sessions, an average of  $M = 21.8$  public posts ( $SD = 24.4$ ) was observed. Given that in almost half of the sessions (48%) no scrolling took place whatsoever, some sessions reveal surprisingly long lists, with up to 174 public posts. Since only public posts were collected, it can be assumed that the actual lists were significantly longer. As expected, the originating sources of the posts varied widely. However, among the top 10 originating

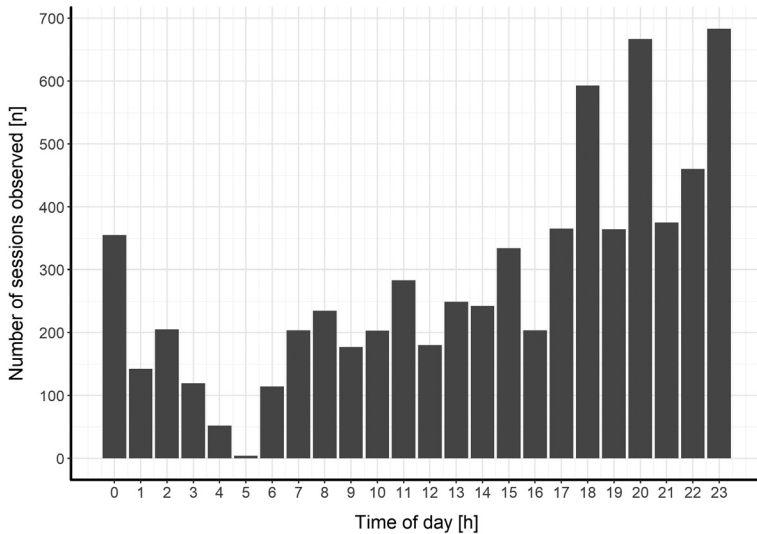


Figure 2. Distribution of observed sessions over the course of a day.

sources for post appearances, four were clearly news outlets (*Stuttgarter Nachrichten*, *Stuttgarter Zeitung*, *tagesschau*, *ZEIT online*), whereas the others represent celebrities (e.g., Neil deGrasse Tyson) or humor/satire (e.g., *ZDF heute-show*). To put this into perspective, all 23 participants reported their Facebook feed in the survey at  $t_1$  as including, on average, 18 per cent news compared to 33 per cent content from various Facebook groups or sites and 40 per cent content from private people; the remaining posts were divided between advertisements and other content (e.g., event announcements). When asked the same question again at  $t_2$  for the previous two weeks (i.e., the observation period), the proportions did not differ substantially.<sup>5</sup>

Interestingly, 6,778 posts were collected more than once within the sample. In other words, although we only observed 23 participants, 68 per cent of all public posts were also shown in other participants' news feeds.

For 337 of the posts, a total of 691 interactions was observed. The large majority of observed interactions involved "liking" public posts: 644 out of the observed 691 interactions fell into this category. The remaining 47 interactions involved clicking a link in the post. This is in line with the forms of online action that respondents claimed to take: 18 out of 23 participants said that they "liked" political posts on Facebook. As 10 out of 23 participants claimed to share political posts and 11 out of 23 respondents claimed to share news posts, future studies should also include computational observation of sharing as an interaction.

**Observational and self-reported data.** When asked about their general Facebook use, 12 of the 23 participants at  $t_1$  and 6 out of the 8 respondents at  $t_2$  claimed to use Facebook multiple times a day. This is in line with the observation findings. It also indicates that heavy users were more inclined to respond as part of the second wave, which was promoted every time that participants navigated to Facebook after the end of the observation period. For most participants, Facebook played a mediocre role as source of political information; that said, for 4 out of the 23 participants, Facebook was indeed the main source of political information. All participants used Facebook notably more often than other social networking sites (i.e., Instagram, Snapchat, Twitter, YouTube).

Although we asked for comparative measures, such as (offline) news outlets to which participants had subscribed and the appearance of posts from these outlets within the observation period, the number of participants did not allow for meaningful analysis of the data we collected. The disappointing response rate, especially at  $t_2$ , also prohibited any analytical

insights into the combination of observation and survey, such as the value of including actually observed content in the questionnaire.

## 5. Discussion

The main goal of this paper was to present an open-source browser plug-in for unobtrusive computational observation of content, context, and online human behavior within algorithmically curated media environments. We discussed the opportunities and challenges regarding computational observation within fragmented and personalized media environments, and we compared methodological approaches that are in current use. Finally, we reported a proof-of-concept study on participants' news exposure on Facebook.

Our plug-in approach has various advantages over other approaches. It allows unobtrusive data collection on a large scale under real-world circumstances. It also enables user interactions to be captured along with the content and the context in which the content was encountered. Moreover, the broad and flexible development of the plug-in does not limit its application to specific websites but allows for combination and comparison across different users, browsers, and settings, and even across different algorithmically curated media environments. Finally, by collecting the content encountered and distributing messages, our approach allows observation to be combined with accompanying surveys and/or content analyses.

However, our plug-in is limited to desktop computers and specific browsers (i.e., Chrome and Firefox). Although it could be deployed to certain other browsers (including Edge and Opera) at almost no additional expense, there remains a concern about the generalizability of this form of computational observation. This primarily includes the absence of Apple's Safari browser. Furthermore, Facebook, in particular, is commonly accessed through smartphone browsers and proprietary apps; as of today, smartphone browsers do not support the easy use of plug-ins, and proprietary apps are not observable with appropriate means. This is especially critical, as actual tracking data from mobile devices has recently revealed even higher overreporting from mobile users than from other study participants (Jürgens, Stark, & Magin, 2019).

It is striking that we faced severe recruitment difficulties. This is in spite of the fact that we employed a previously successful survey panel (twice), informed potential participants about all aspects of the project, provided proportionate incentives, made plug-in installation as easy as possible

through plug-in stores, and provided participants with complete sovereignty over their data. Although a handful of potential participants reported that they mainly used Apple Safari for browsing the web, we are unable to explain more generally whether the low response rates were due to technological issues or to other considerations, such as privacy concerns. In addition to the very low response rates, we also faced the mortality rates that are common within panel studies. Similar endeavors in future should keep this in mind; researchers might try applying stronger incentives, such as (financial) allowances, as well as developing a plug-in version for Apple Safari.

As outlined above, the use of computational observation requires consideration of a range of technological, ethical, and practical aspects. Importantly, using a browser plug-in to collect a large amount of potentially sensitive data requires careful consideration of data security and privacy issues. We addressed these issues by complying with the European Union's GDPR during all phases of the project. However, data collection also needs to comply with the terms of use of specific algorithmically curated environments and to take participants' concerns about data privacy seriously. Participants should therefore be given constant access to the data collected from them, and they should have the opportunity to delete data at any point in time.

By following these suggestions, results from computational observation are likely to provide in-depth insights, as our proof-of-concept study indicates. Given the exploratory character of this study, our findings are preliminary and should be handled with adequate care. The study, however, examines the applicability and validity of such a research setup, and a variety of contemporary research questions may benefit from this methodological approach. For example, media use studies might increase data reliability by employing observational data rather than solely self-reported data. This approach might also help in identifying content, which participants sometimes prefer not to report on (e.g., politically extreme content). Studies regarding filter-bubble or echo-chamber phenomena may identify and examine mechanisms of algorithmically curated media environments by using a combination of computational observation, panel survey, and (automated) content analysis. Also, laboratory-based observational studies, such as on selective exposure, could be freed from their dependency on laboratories and rather build on computational observation to be conducted on the participants' own devices.

Taken together, as the plug-in can be adapted to any website, computational observation can easily be employed for essentially any empirical study where it is necessary to observe participants' online behavior, within

or outside of social networking sites. This potential methodological generalizability not only allows to cut on development resources, but it also enables future projects to collect data across platforms—a research demand that has been raised repeatedly (e.g., van Atteveldt & Peng, 2018; Wallach, 2016).

## 6. Author Note

The authors would like to thank the two anonymous reviewers and the editor for their suggestions, which have significantly improved the manuscript. The authors would also like to thank Phelia Weiß, who helped to conduct the proof-of-concept study, and Kai Hilgers, who helped to develop the plug-in.

## Notes

- 1 <https://github.com/MarHai/fbforschung>
- 2 [https://github.com/MarHai/fbforschung/tree/master/json\\_schemas/](https://github.com/MarHai/fbforschung/tree/master/json_schemas/)
- 3 <https://chrome.google.com/webstore/detail/fbforschungde-browser-plu/faemgdmnkflbiakkkchdgpaphljccpch>
- 4 <https://addons.mozilla.org/firefox/addon/fbforschung/>
- 5 Statistical inferences are difficult, because only 8 cases included data for both  $t_1$  and  $t_2$ .

## 7. References

- Ahmadi, M., & Wohn, D. Y. (2018). The antecedents of incidental news exposure on social media. *Social Media + Society*, 4(2), 1–8. <https://doi.org/10.1177/2056305118772827>
- Anderson, A. A., Brossard, D., Scheufele, D. A., Xenos, M. A., & Ladwig, P. (2014). The “nasty effect”: Online incivility and risk perceptions of emerging technologies. *Journal of Computer-Mediated Communication*, 19(3), 373–387. <https://doi.org/10.1111/jcc4.12009>
- Bodó, B., Helberger, N., Irion, K., Zuiderveen Borgesius, F. J., Möller, J., van de Velde, B., ... de Vreese, C. H. (2018). Tackling the algorithmic control crisis - the technical, legal, and ethical challenges of research into algorithmic agents. *Yale Journal of Law and Technology*, 19(1), 133–180.
- Bozdag, E. (2013). Bias in algorithmic filtering and personalization. *Ethics and Information Technology*, 15(3), 209–227. <https://doi.org/10.1007/s10676-013-9321-6>
- Diakopoulos, N. (2014). Algorithmic accountability. Journalistic investigation of computational power structures. *Digital Journalism*, 3(3), 398–415. <https://doi.org/10.1080/21670811.2014.976411>
- Eslami, M., Rickman, A., Vaccaro, K., Aleyasen, A., Vuong, A., Karahalios, K., ... Sandvig, C. (2015). “I always assumed that i wasn’t really that close to [her]”: Reasoning about invisible algorithms

- in news feeds. *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems - CHI '15*, 153–162. <https://doi.org/10.1145/2702123.2702556>
- European Parliament, & Council of the European Union. *Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)*, (2016).
- Freelon, D. (2018). Computational research in the post-API age. *Political Communication*, 35(4), 665–668. <https://doi.org/10.1080/10584609.2018.1477506>
- Gedeshi, I., Zulehner, P. M., Faradov, T., Rotman, D. G., Swyngedouw, Marc (Flanders), Fotev, G., ... European Values Study Group. (2010). *European Values Study 2008: Integrated dataset (EVS 2008)* [Data set]. <https://doi.org/10.4232/1.10059>
- Guess, A. M. (2015). Measure for measure: An experimental test of online political media exposure. *Political Analysis*, 23(1), 59–75. <https://doi.org/10.1093/pan/mpu010>
- Guess, A. M., Munger, K., Nagler, J., & Tucker, J. (2018). How accurate are survey responses on social media and politics? *Political Communication*, 0(0), 1–18. <https://doi.org/10.1080/10584609.2018.1504840>
- Haenschen, K. (2019). Self-reported versus digitally recorded: Measuring political activity on Facebook. *Social Science Computer Review*. <https://doi.org/10.1177/0894439318813586>
- Haim, M., Graefe, A., & Brosius, H.-B. (2016). *The burst of the bubble? Effects of automated personalization on news diversity*. Paper presented at the 65th meeting of the International Communication Association presented at the Fukuoka. Fukuoka.
- Haim, M., Kümpel, A. S., & Brosius, H.-B. (2018). Popularity cues in online media: A review of conceptualizations, operationalizations, and general effects. *Studies in Communication and Media*, 7(2), 186–207. <https://doi.org/10.5771/2192-4007-2018-2-58>
- Hasebrink, U., & Popp, J. (2006). Media repertoires as a result of selective media use. A conceptual approach to the analysis of patterns of exposure. *Communications*, 31(3), 369–387. <https://doi.org/10.1515/commun.2006.023>
- Jürgens, P., Stark, B., & Magin, M. (2019). Two half-truths make a whole? On bias in self-reports and tracking data. *Social Science Computer Review*. <https://doi.org/10.1177/0894439319831643>
- Kalogeropoulos, A., & Newman, N. (2017). *“I saw the news on Facebook”: Brand attribution when accessing news from distributed environments*. Oxford: Reuters Institute for the Study of Journalism.
- Kim, S. J. (2014). A repertoire approach to cross-platform media use behavior. *New Media & Society*, 18(3), 353–372. <https://doi.org/10.1177/1461444814543162>
- Kümpel, A. S. (2018). The issue takes it all? Incidental news exposure and news engagement on Facebook. *Digital Journalism*, 5(1), 1–22. <https://doi.org/10.1080/21670811.2018.1465831>
- Leiner, D. (2014). *Convenience samples from online respondent pools: A case study of the SoSci Panel*. Retrieved from <https://www.researchgate.net/publication/259669050>
- Menchen-Trevino, E. (2016). Web Historian: Enabling multi-method and independent research with real-world web browsing history data. *IConference 2016 Proceedings*. Presented at the iConference 2016, Philadelphia, PA. <https://doi.org/10.9776/16611>
- Microsoft. (2017, August 2). Porting an extension from Chrome to Microsoft Edge. Retrieved November 13, 2018, from <https://docs.microsoft.com/en-us/microsoft-edge/extensions/guides/porting-chrome-extensions>
- Möller, J., van de Velde, R. N., Merten, L., & Puschmann, C. (2019). Explaining online news engagement based on browsing behavior: Creatures of habit? *Social Science Computer Review*. <https://doi.org/10.1177/0894439319828012>

- Newman, N., Fletcher, R., Kalogeropoulos, A., Levy, D. A. L., & Nielsen, R. K. (2018). *Digital news report 2018*. Retrieved from <http://media.digitalnewsreport.org/wp-content/uploads/2018/06/digital-news-report-2018.pdf>
- Newman, N., Fletcher, R., Levy, D. A. L., & Nielsen, R. K. (2016). *Digital news report 2016*. Retrieved from <http://reutersinstitute.politics.ox.ac.uk/sites/default/files/Digital-News-Report-2016.pdf>
- Price, V., & Zaller, J. (1993). Who gets the news? Alternative measures of news reception and their implications for research. *Public Opinion Quarterly*, 57(2), 133–164. <https://doi.org/10.1086/269363>
- Prior, M. (2009). The immensely inflated news audience: Assessing bias in self-reported news exposure. *Public Opinion Quarterly*, 73(1), 130–143. <https://doi.org/10.1093/poq/nfp002>
- Reinecke, L., & Hofmann, W. (2016). Slacking off or winding down? An experience sampling study on the drivers and consequences of media use for recovery versus procrastination. *Human Communication Research*, 42(3), 441–461. <https://doi.org/10.1111/hcre.12082>
- Revilla, M., Ochoa, C., & Loewe, G. (2017). Using passive data from a meter to complement survey data in order to study online behavior. *Social Science Computer Review*, 35(4), 521–536. <https://doi.org/10.1177/0894439316638457>
- Roßteutscher, S., Schmitt-Beck, R., Schoen, H., Weßels, B., Wolf, C., Preißinger, M., ... Gärtner, L. (2018). *Wahlkampf-Panel 2017 (GLES)* [Data set]. <https://doi.org/10.4232/1.13150>
- Sax, L. J., Gilmartin, S. K., & Bryant, A. N. (2003). Assessing response rates and nonresponse bias in web and paper surveys. *Research in Higher Education*, 44(4), 409–432. <https://doi.org/10.1023/A:1024232915870>
- Scharnow, M. (2016). The accuracy of self-reported internet use—a validation study using client log data. *Communication Methods and Measures*, 10(1), 13–27. <https://doi.org/10.1080/19312458.2015.1118446>
- Struckmann, S., & Karnowski, V. (2016). News consumption in a changing media ecology: An MESM-study on mobile news. *Telematics and Informatics*, 33(2), 309–319. <https://doi.org/10.1016/j.tele.2015.08.012>
- Taneja, H., Webster, J. G., Malthouse, E. C., & Ksiazek, T. B. (2012). Media consumption across platforms: Identifying user-defined repertoires. *New Media & Society*, 14(6), 951–968. <https://doi.org/10.1177/1461444811436146>
- van Aelst, P., Strömbäck, J., Aalberg, T., Esser, F., de Vreese, C. H., 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. <https://doi.org/10.1080/23808985.2017.1288551>
- van Atteveldt, W., & Peng, T.-Q. (2018). When Communication Meets Computation: Opportunities, Challenges, and Pitfalls in Computational Communication Science. *Communication Methods and Measures*, 12(2–3), 81–92. <https://doi.org/10.1080/19312458.2018.1458084>
- Wallach, H. (2016). Computational social science: Toward a collaborative future. In R. M. Alvarez (Ed.), *Computational Social Science* (pp. 307–316). <https://doi.org/10.1017/CBO9781316257340.014>
- Webster, J. G., & Ksiazek, T. B. (2012). The dynamics of audience fragmentation: Public attention in an age of digital media. *Journal of Communication*, 62(1), 39–56. <https://doi.org/10.1111/j.1460-2466.2011.01616.x>
- Wolfsfeld, G., Yarchi, M., & Samuel-Azran, T. (2016). Political information repertoires and political participation. *New Media & Society*, 18(9), 2096–2115. <https://doi.org/10.1177/1461444815580413>
- Zuiderveen Borgesius, F. J., Trilling, D., Möller, J., Balázs, B., de Vreese, C. H., & Helberger, N. (2016). Should we worry about filter bubbles? *Internet Policy Review*, 5(1), 1–16. <https://doi.org/10.14763/2016.1.401>

## About the authors

**Mario Haim:** Department of Media and Social Sciences, University of Stavanger, Norway

Correspondence address: [mario.haim@uis.no](mailto:mario.haim@uis.no)

**Angela Nienierza:** Department of Media and Communication, LMU Munich, Germany

Creative Commons License CC BY

(<https://creativecommons.org/licenses/by/4.0/>)