# PRACTICAL AI TRANSPARENCY: REVEALING DATAFICATION AND ALGORITHMIC IDENTITIES

Ana Pop Stefanija & Jo Pierson*

**ABSTRACT**

How does one do research on algorithms and their outputs when confronted with the inherent algorithmic opacity and black box-ness as well as with the limitations of API-based research and the data access gaps imposed by platforms' gate-keeping practices? This article outlines the methodological steps we undertook to manoeuvre around the above-mentioned obstacles. It is a "byproduct" of our investigation into datafication and the way how algorithmic identities are being produced for personalisation, ad delivery and recommendation. Following Paßmann and Boersma's (2017) suggestion for pursuing "practical transparency" and focusing on particular actors, we experiment with different avenues of research. We develop and employ an approach of *letting the platforms speak* and *making the platforms speak*. In doing so, we also use non-traditional research tools, such as transparency and regulatory tools, and repurpose them as *objects of/for study*. Empirically testing the applicability of this integrated approach, we elaborate on the possibilities it offers for the study of algorithmic systems, while being aware and cognizant of its limitations and shortcomings.

Keywords: datafication; algorithmic identity; practical transparency; methodology; digital methods; subject access request.

* imec-SMIT, Vrije Universiteit Brussel, Belgium.

## 1    INTRODUCTION

Today, there is almost no area in everyday life that has not been mediated or impacted by Artificial Intelligence (AI). From recommender systems for news, apps, routes, products, to job applications, financial services, health care, education, criminal justice, etc., individuals have been increasingly, to lesser or greater degree, subjected to the automated decision making (ADM) by some kind of algorithmic and AI systems. More and more decisions impacting individuals are based on what we can call 'algorithmic identity' (Cheney-Lippold, 2011; but also Jarrett, 2014; Reigeluth, 2014) — guided by extensive profiles about people, uncovering their affinities and interests and predicting their behaviour. With the ubiquity of these ADM and AI systems, it becomes an issue of urgency to be able to investigate them, reveal their workings, and explain their outputs and impact.

We could say that this article is a "byproduct" of our attempt to investigate how algorithmic identities are being produced by few sampled actors (*Facebook, Google, Quantcast* and *Oracle)* participating in the process of algorithmic identity building for personalisation, ad delivery and recommendation. For us that meant a critical investigation of the inner workings of algorithmic systems, of the datafication practices that enable the algorithmic identity creation, in particular the actors participating in these processes, the types of data that are used, the sources of data and — importantly — their relation to the inference-making processes that are building blocks for the algorithmic identities. But an analytical inquiry like this confronted us with the question of how to investigate these issues? How does one do research on algorithms and their outputs when confronted with the inherent algorithmic opacity and black box-ness as well as with the limitations of API-based research and the data access gaps imposed by platforms' gate-keeping practices? How does one overcome and manoeuvre the limitations when dealing with data provided or extracted from these platforms while being aware of and critical of the 'methodological bias' (Marres and Gerlitz, 2015), the un-rawness of data (Gitelman, 2013) and the level of mediation (van Es et al., 2013)? This had led to our research focus of investigating the process of how an individual algorithmic is being created by few different platforms. To achieve this, we experiment with novel avenues for methodologically investigating this process and the underlying datafication processes and practices. In doing so, we also contribute towards answering the question of how to uncover and explain datafication and algorithmic identity.

This paper documents and discusses our attempts and experimentations in studying and researching algorithms while doing digital social research (Lindgren, 2019). We rely on data collection methods that manoeuvre around the API restrictions and make use of non-

traditional data sources, like transparency and regulatory tools. We do this by using a mixed method design for which we developed and adopted two approaches: *letting the platforms speak* and *making the platforms speak*. This led to an investigation on two levels, interface and software, while employing two corresponding overarching methods, technography and digital methods (Rogers, 2017). Experimenting with this methodological setup, our experience and results show that, while there are limitations, this approach enables an in-depth critical inquiry and generates valuable new insights into the processes of algorithmic construction of identity, data extraction and inferential analytics.

What follows is an outline of the techniques we employed around the limitations we encountered when dealing with platforms' algorithmic systems for the purpose of research. We see this approach as just one possible path for doing research *on* algorithms, AI and platforms. As such, it does not aim to be taken as a generalizable, "apply-to-all" approach, but it aims foremost to inspire, test and experiment, and explore the possibilities and limits of different approaches and tools. We will elaborate on the rationale behind the approach, the methods chosen, the particular tools and research protocols applied, as well as the specific steps taken. First, we discuss some recent developments in digital social research, the restrictions imposed by platforms and the very nature of algorithms, and elaborate how that impacts the ability to do digital social research. This is followed by outlining our methodological design and rationale and detailing the particular steps and tools we used. A substantial part is dedicated to the elaboration of our results. We conclude with a discussion on the advantages and limitations of the particular methodological choices and tools.

## 2    PRACTICAL AI TRANSPARENCY

The networked infrastructure of the internet, with its technological capacity to track user movements across different web sites, apps and servers, has given rise to an industry of web analytics firms that are actively amassing information on individuals and fine-tuning computer algorithms to make sense of that data. Via the process of *datafication* – 'the transformation of the social actions of their users to quantified data' (Mayer-Schönberger and Cukier, 2013, p. 78) – and the collection of data via tracking technologies, combined with the analytics capabilities of algorithms and companies, the aim of many of these companies it to create what Cheney-Lippold (2011) calls 'algorithmic identity' – "an identity formation that works through mathematical algorithms to infer categories of identity on otherwise anonymous beings" (p.165). Datafication can be understood both as *action* and as *aim*. As an action, it means "transformation of social action into

online quantified data, thus allowing for real-time tracking and predictive analysis" (van Dijck, 2014, p.198). As an aim, it relates to the pursuit to collect, monitor, analyse, understand and use people's behaviour for behaviour prediction, affinity profiling, but also for 'unstated preset purposes' (van Dijck, 2014, p.205). Raley (2013) calls the latter 'data speculation', i.e. a value yet to be added to the data and 'informational patterns still to come' (p. 123). This is closely tied, first, with the belief "in the objective quantification and potential tracking of all kinds of human behaviour and sociality through online media technologies" (ibid., p.198) to which van Dijck refers to as 'dataism', and second, with the 'collect everything' approach (van Dijck, 2014; Sadowski, 2019; Andrejevic & Gates, 2014). The creation of an algorithmic identity is possible because of the process of datafication, and datafication and dataism are the building blocks for behaviour prediction and affinity profiling, among other things, for targeted advertising and personalisation. However, this algorithmic identity is a construct, "it is not the personal identity of the embodied individual but rather the actuarial or categorical profile of the collective which is of foremost concern' to new, unenclosed surveillance networks" (Hier, 2003: 402 in Cheney-Lippold, 2011, p. 177). So how do we investigate the process of algorithmic identity and the underlying processes of datafication?

In his recent article Axel Bruns (2019) (rightfully) states that the APIcalypse has arrived and it seriously impacts our ability as social science researchers to critically study society via the digital. This results in a restriction as regards who has access to the platforms' data via their Application Programming Interfaces (APIs), which makes access to data either impossible or possible only for the chosen few and under strict conditions. Hence the APIcalypse limits the possibilities to inspect and investigate phenomena happening "in the digital". The importance of this gatekeeping is even greater if we consider that the online is never a separate realm, as decisions made about individuals based on the digital traces they leave behind can impact their offline lives as well. Being severely restricted and limited in investigating the digital and the algorithmic, seriously impacts the ability of researchers and scientists to investigate and criticise these systems, hold to account their proprietors and remedy their outputs.

To borrow the definition by Venturini and Rogers (2019), API based research is

> an approach to computational social sciences and digital sociology based on the extraction of records from the datasets made available by online platforms through their application programming interfaces (or APIs). This type of research has allowed the collection of detailed information on large populations, thereby effectively challenging the distinction between qualitative and quantitative methods. (p. 1).

As such, this approach enabled the studying of a variety of phenomena pertaining to the interplay and mutual influence of both technology and society, and mediated numerous findings, previously not possible at such a large scale. However, this is not the only approach undertaken or proposed by researchers and scholars for studying the digital, or the best one. As Venturini (2018) notes, 'when all you have is a Twitter feed, everything looks like a hashtag' (p. 4210). This refers to the limitations imposed by the affordances of the platforms when we see them as objects of research. We use these statements as an entry point to discuss some of the approaches developed for studying the relationship between the digital, the algorithmic and the societal, their limitations and shortcomings. In the paragraphs that follow we briefly outline some of them, and outline our own developed methodological approach, as being a response to both the APIcalypse and the dominant discourse of API-dependability for research.

The approach of *auditing algorithms* was proposed by Sandvig et al. (2014) and entails different techniques to uncover the inner workings of algorithmic agents. Depending on the infrastructure and affordance of the system, the objective (input, output or system) and available resources, these techniques range from relying on APIs, use of software and hardware infrastructures to users' input, either to investigate the code or the outputs of the system. Weltevrede (2016) talks about the adoption of a *device-driven approach*, as a way to focus on the 'the specific strategies or intents embedded in algorithms' (p. 106) and to repurpose the 'analytical affordances of the algorithmic systems/devices' (ibid.). Because algorithms are techno-epistemological devices, the analytical inquiry is dependent on the system's affordances, so on what the system allows and limits to be seen. As such it requires a combination of different types of methodological and conceptual resources to study device-captured data points. This approach shares similarities with the *reverse-engineering* one, as a diagnostic approach that allows for an observation of the relationship between the inputs and outputs, and a way to obtain 'missing knowledge' (Bucher, 2012, p. 79) and grasp a model of how the particular algorithmic system works. As a strategy to see to *what* the algorithm pays attention to, is a "process of articulating the specifications of a system through a rigorous examination drawing on domain knowledge, observation, and deduction to unearth a model of how that system works." (Diakopoulos, 2014, p. 13).

All these approaches are characterised by a move away from the quest to open the black box towards investigating *algorithms in action*, at work, in practice. It is a quest towards 'unknowing algorithms' (Bucher, 2018; Annany and Crawford, 2016), studying them as 'part of specific situations' (Bucher, 2018, p. 49) and uncovering the actor-network assemblages/configurations (Annany and Crawford, 2016). By observing the effects of the system, researchers are able to overcome the obstacles of

the "black box", and to assess the 'operational principles of systems' (Bucher, 2012, p. 77) and its actual working. Additionally, investigating algorithms as an assemblage(s), to borrow Annany & Crawford's (2016) suggestion, is to look at them *as a system* and *across a/the system*.

However, this doesn't solve (all of) the difficulties of socially investigating algorithms. Algorithms are predominantly patent protected and proprietary software, with their inherent opacity stemming from the underlying machine-learning process at work. It is never a single algorithm, but always an algorithmic *system* of interconnected and interrelated algorithms (Gillespie, 2014; Bucher, 2018). In addition, these systems are in a 'perpetual beta' state (Weltevrede, 2016) with constant and continuous A/B testing, fine tuning and upgrades, making the study of algorithmic systems almost a study of a 'historical object' (Bucher, 2012). All this coupled with the research affordances (Rogers, 2013) of — and the restricted access to — the platforms' APIs and code crucially limits and impacts digital social research and complicates the task of developing and applying the appropriate methodological apparatus and tools.

Faced with these inherent characteristics, the calls for transparency of algorithmic systems, initially aiming towards total transparency, have shifted their focus significantly. Paßmann and Boersma (2017), differentiate between two notions of transparency. *Formalised transparency*, which would like to see more inside the content of the black box and 'obtain more positive knowledge' (p.140) and *practical transparency*, which does not try to open the black box, but to 'develop skills without raising the issue of openability' (ibid.). These skills should help researchers deal with the (parts of the) algorithms that we still don't have knowledge about, and probably we won't be able to have. Thus, the aim is actually to ask and investigate how to 'behave towards what remains black after all.' (ibid., p. 140). In order to find ways to work around these unknowns, the authors suggest other sources, external to the algorithms that will help turn 'unknown unknowns in to known unknowns' (p. 145), such as ethnographic data or other sources that are some kind of everyday knowledge. Our research follows the principle of practical transparency.

## 3 METHODOLOGICAL PATHWAYS THROUGH THE (ALGORITHMIC) SYSTEM

In his book *Design research and the new learning*, Buchanan (2001) states:

> By definition, a system is the totality of all that is contained, has been contained, and may yet be contained within it. We can never see or experience this totality. We can only experience our personal pathway through a system. (p. 12).

This corresponds with the methodological and sampling approach that we adopt in our empirical research: zooming in on a few platforms but looking at the wider system/assemblages of actors participating in the creation of an algorithmic identity of a single individual. Focusing on a research subject of one, we also expand our research to the other social and technological actors partaking in the process. In this following section we elaborate on the methodological setup while discussing the specific aspects we took into consideration and the limitations and opportunities we were faced with.

*Methodological design.* Our methodological approach is the result of a two-way process. First, we built our research on an assessment of the analytical affordances of the platforms in our study and of the mechanisms and tools known and available to the researchers. We tested and experimented with a variety of digital methods and tools, ranging from API-access ones to scraping ones. Through a process of going back and forth, we finalised the list of tools based on their applicability to the research questions and their particular affordances, while constantly being aware of their level of mediation (van Es et al., 2018).

Next, we experimented with the method of an *interface walkthrough* — where we mimicked and rehearsed ordinary use (*researchers as users* perspective) (Dieter and Tkacz, 2020). In that way we investigated what could be collected and used as data for research through what was available via the interface of the platforms. However, if we were to experiment with "out of the (black) box" approaches and tools, we had to think both more critically and creatively. In doing so, we "took advantage" of the newly established regulatory and transparency mechanisms and repurposed them as objects/tools for study. The platforms we queried have developed and made available (confined) gateways to transparency and explainability, as an attempt to provide more information on data collection and personalisation practices. We decided to experiment with these transparency and accountability tools and see if we could repurpose them as objects for study. Additionally, we were curious to investigate how the General Data Protection Regulation's (GDPR) Article 15, its' corresponding recitals and in particular the Data Subject Access Request mechanism (European Commission, 2016) could be used for academic research.

*Approach.* Faced with the above-mentioned challenges, one of the strategies created, tested and employed was to work with what is available and be creative in finding ways to do research relying on the affordances of platforms themselves and repurposing transparency and regulatory tools as objects for/of study.[1] We define these approaches as *letting the platforms*

---

[1] Data was collected from *Facebook*, *Google*, *Oracle*, *Quantcast* and visited webpages. The automated tools used were *AdAnalyst*, *TrackerObserver*, *PriBot* and *Privacy Score*. Data recorded included capturing of trackers on websites, *Facebook* Ads Shown, *Facebook*

*speak* and *making the platforms speak,* focusing on achieving *practical transparency* (Paßmann and Boersma, 2017) through investigating *algorithms in action* and studying them by observing their outputs and effects.

*Sampling.* The insights collected and discussed in this paper are the result of the data originating from one research subject (n=1). It is collected via different means over a period of six months[2]. Choosing the personalised, one-research-subject-only approach, allows for the observation of real user-algorithmic agents interactions, where "pre-existing profiles, browsing histories, technology fingerprints, and other organically developed profile information are used." (Bodo et al., 2018, p. 143). This real-world observation is advantageous in comparison with the use of sock puppet audits or dummy users, as it overcomes the shortcomings of 'non-adequate approximations of real-life users' (ibid.), allowing for investigation of the effects of algorithmic agents on individual users (ibid., p. 144). As such, the detailed (data) account of a personalised experience offers an overview of 'the whole spectrum of online and offline, personalised and non-personalised information flows.' (ibid., p. 145). Additionally, the insights offered by small data bear the quality of more context-aware research, granularity and depth of the data and the findings by combining various methods, complementing data and triangulating the findings (Crawford, 2013). As the method and type of data should follow the research question (Van Es et al., 2017), small data gathered using digital data analysis enables for a qualitative and contextualised investigation (Kitchin, 2014).

Focusing on algorithms in actions, around a user (a real individual, with a browsing history and data scattered around the digital space and different online and offline databases), that exhibits real-life behaviour and for whom information "in the wild" already exists, enables not just for non-lab experimentation, but also for fully taking advantage of additional non-traditional research tools, such as Data Subject Access Requests. We are aware that one of the difficulties with the auto-technographic approach is its highly individualised and personalised approach "as the observation of the interface is confined to the 'me-centric view of the researcher's own

---

Interests, Interactions with Advertisers (*Facebook*), Advertisers that have uploaded contact details (*Facebook*), Why am I seeing this Ad (*Facebook*), assigned interests by *Google* and reasons for assigning them. Data was collected in the period of November 27, 2018 to June 6, 2019, with different recording periods for different insights, following the browsing behaviour of one research subject. A research web browser was set prior to the start of the data collection process.

[2] Data was collected in the period from November 27, 2018 to June 6, 2019, from Brussels, Belgium. The data collection period however differs between the different tools used and the related observation. This is elaborated in more detail in the sections related to each particular tool.

account" (Weltevrede, 2016, p. 107). However, this approach both enables to manoeuvre around the "black-boxed" systems and to follow the advice by boyd and Crawford (2012) that 'the size of data should fit the research question being asked' (p. 670).

## 4    RESEARCHING ALGORITHMS IN ACTION

*Letting the platform speak* approach relies on what the platforms themselves allow to be seen and to be visible at an interface level, without the assistance or help of additional data collection tools, relying on the affordances of the presentation layer of the platforms and their *front-end*. Literary it means looking at what information platforms willingly provide and reveal via the user interface. This approach also helps uncover the platform's *politics of visibility*, i.e. what the algorithmic system itself decides to make visible and the insights it permits willingly. In addition to the focus on the interface, this approach entails use of external available sources that describe and reveal the workings of the system (Bucher, 2012a, p. 74): technical documentation, specifications, patents, media talks, but also help sections for users and advertisers. However, what we did in a novel way, and where we add to the repository of methods for research is the usage of the *transparency tools* enabled by platforms (such as Ad Settings, Data Explanations and Ad Explanations), the *privacy policies* and the *Data subject Access requests*, enabled by the GDPR. In that sense, we employed a 'multi-site technography' (Bucher, 2012, p. 73): as algorithmic systems are always assemblages and always in interaction with other actors and systems, be it technical or human, all these "sites" can be used as sources of data and insights for digital social research.

Data collection-methods wise, technography, as defined by Taina Bucher (2012) was adopted, as "a descriptive-interpretative approach to the understanding of software, rooted in a critical reading of the mechanisms and operational logic of technology." (p. 71). This is employed via observation, where the daily changes of the information provided by the platforms are observed and recorded. This approach was chosen because it allows for a granular, detailed dossier of the interaction and communication between the user and the system, it enables for insights into the actors they are *in communication* with, into what is 'the interplay between a diverse set of actors (both human and nonhuman)' (Bucher, 2012, p. 69). This is especially important for the investigation of the actor-network around the data collected and sources used for affinity profiling and algorithmic identity-building, their position within the network and in 'particular sociotechnical events' (Latour, 2005: 128 in Bucher, 2012, p. 72).

The *making the platform speak* approach, on the other hand, looks for insights not relying on what the software makes visible willingly, but by

*forcing* the software to reveal itself and its inner workings. It relies on the use of automated scraping and crawling tools and tools relying on platforms' Application Programming Interfaces (APIs). In that sense it also makes visible the *politics of knowledge* of the platform, i.e. what the platform allows to be known, if one has the knowledge and tools to seek knowledge. While this can be more insightful, it is still limited. This approach aims to make the system reveal itself, in order to gain more in-depth knowledge or insights by looking not just how it produces outputs, but also to uncover things not visible at an interface level and to the human eye. In that regard, this is an analysis done at a *software* level. This approach implies that the algorithmic devices and systems will be *forced* to speak, meaning, the "analytical gaze" goes beneath the surface and what is visible and tries to uncover some inner workings of these systems.

We specifically set up a research browser through which the platforms and other actors would be able to gather as much possible information on the behaviour, actions, patterns of behaviour of the research subject and thus provide personalised search results, ads and recommendations. This enabled us to – as objectively as possible – investigate the datafication practices and the creation of algorithmic identity, while being aware of the multitude of factors affecting data collection and algorithmic outputs in the form of personalisation. In addition we were able to further investigate the assigned algorithmic identity via the outputs provided both by the used search engine (*Google)* and browser (*Chrome*) and the platforms visited during the period of the data collection phase.[3] Steps were also undertaken to allow for as much data collection and data sharing between *Facebook* and third-parties as possible, by setting up the preferences, permission and settings options[4].

---

[3] We set up the research browser by installing a "clean browser", deleting all the previous cookies, browsing history and preferences, and setting the preferences to allow for a maximum data collection by the platform and associated third parties: cookies were enabled, keeping record of web and app activity and location was enabled (location history, device information - info about contacts, calendars, apps, and other device data to improve users' experience across *Google* services, voice and audio activity, *YouTube* search History, *YouTube* watch history), as well as "Chrome browsing history and activity from websites and apps that use Google services" (that includes: activity from sites and apps that partner with *Google* to show ads; *Chrome* history (if Chrome Sync is turned on; app activity, including data that apps share with *Google*; Android usage & diagnostics, like battery level, how often you use your device and apps, and system errors). Ad settings were adjusted too, enabling ad personalisation, giving *Google* permission to show ads based on user's activity on *Google* services (such as Search or *YouTube*) and websites and apps that partner with *Google* to show ads. Whenever a consent by websites was asked in regard to data collection (in accordance with the GDPR), consent was given.

[4] The steps we took to set up and allow *Facebook* to maximise the data collection for the research subject were the following: changing the privacy settings and enabling data

# 5   INVESTIGATING DATAFICATION AND ALGORITHMIC IDENTITIES

We start our analysis by investigating datafication practices and the network of actors around the research subject. This is an important starting point, as the creation of an algorithmic identity relies on behavioural data collected via tracking elements present on both the web[5] and in apps. This step, additionally, guides the further analysis of the process of algorithmic identity creation: what data is seen as a worthy signal and what behaviour is taken as important/proxy for affinity profiling – 'grouping people according to their assumed interests rather than their personal traits' (Wachter, 2019, p. 33), based on proxies (friends, likes, groups, IP address and similar). Importantly, we are interested to see if only 'raw' data is taken as basis for inferences or there are other (hidden) mechanisms and 'cooked' data (Gitelman, 2013). The structuring of the results follows the same path: we first elaborate on our approach and findings regarding datafication and then focus on methods to investigate and assess algorithmic identity.

## 5.1   Investigating datafication

In order to investigate the formation of an algorithmic identity, our first step was to investigate the datafication practices surrounding a user. This provided us with insights into two interrelated aspects: the sources taken as input for the prediction outputs – the 'qualities, preferences, characteristics, intentions, needs and wants of users' (Lehtiniemi, 2016, p.4), affinities and interests — and the network of companies that collect (behavioural) data about the user (traces of user actions and interactions), as well as their dominance and variety. For this we used diverse sources of insight, collecting data on different levels (interface and software) and using a mixed method approach. We did this according to the following consecutive phases: first, using automated tools to record tracking behaviour and data collection, after which, we used privacy policies as source of information regarding data collection practices of platforms and companies. Lastly, we used transparency and regulatory tools as objects for studying datafication practices.

---

collection and data-sharing between Facebook family of companies and services; allowed "Ads based on data from partners"; "Ads based on your activity on Facebook Company Products that you see elsewhere"; allowed Facebook Audience Network; with the setup of the Research browser to enable third-party tracking, *Facebook* was granted access to the full browsing, off-*Facebook*, behaviour of the research subject.

[5] In this research we focus only on tracking datafication actors via web platforms.

| SOFTWARE LEVEL | | INTERFACE LEVEL (repurposing tools as objects of study) | |
|---|---|---|---|
| TrackerObserver | most prevalent companies | Transparency tools (Facebook & Google) | Type of data collected for affinity profiling |
| Privacy Score | Important actors (data brokers) | Privacy Policies | Sources of data and mechanisms of data-sharing |
| | | Data Subject Access Right (Oracle, Quantcast) | Actors with personal data about the user |

| Tools used \| Insights gathered | Tools used \| Insights gathered |

*Figure 1. Overview of the tools used for investigating datafication and insights gathered*

Firstly, by using digital methods and tools, we collected information on the third-party trackers using the browser extension *TrackingObserver*[6] and the automated web scanner *Privacy Score*[7]. This was done at a software level. Both tools offer different insights in correspondence with their aim, affordances and information structure. As a result, they are suitable for different aspects and levels of analysis. Because of the ability to track every browsing behaviour around a particular user, *TrackingObserver* enables investigation of the network of third-party trackers and companies around a particular user and their unique browsing behaviour. From the data collected during a six months period[8], we obtained valuable insights into the network relations and data exchange practices of a multitude of actors. The latter was later used as a source for further investigation.

We triangulated the data obtained via the initial data collection with data available from other sources (*WhoTracksMe*[9] and *Better.fyi*[10]), providing us with several valuable insights: it enabled us to reveal the companies behind the trackers and analyse their presence, to detect the type of trackers

---

[6] Information about the tool is available at: https://trackingobserver.cs.washington.edu/. Last accessed January 29, 2020.

[7] Information about the tool is available at: https://privacyscore.org/. Last accessed January 29, 2020.

[8] Data was collected in the period November 27,2018 – June 4, 2019, and the analysis showed the presence of 4,691 tracking instances observed, set on 287 websites (on average 16,3 trackers per website), with 1,067 unique tracking domains.

[9] We were specifically looking at the Trackers analysis (tracking type and tracker category) and the companies indicated as owning the particular trackers. Information can be found following this link: https://whotracks.me/trackers.html.

[10] We were interested and recording the particular type of trackers, as well as the company owning the trackers. Information can be found following this link: https://better.fyi/trackers/

and their particular purpose. The analysis showed the dominance of a few companies in the network, representing the majority of trackers on the visited websites (Figure 2).

However, we also observed a long tail of many different actors (a large number of trackers with low websites frequency) that captured data about the user's behaviour, supporting similar findings by Binns et al. (2018b). Categorising the detected trackers based on a taxonomy, we discovered a presence of a vast and well-developed network of *ad networks*, counting for more than half (57.23%) of the detected unique trackers. These findings are important for several reasons. The detected long tail is worrying as it indicates that a great number of companies get some and partial data from the research subject and users in general. This is even more of a cause for concern as the profiling-oriented businesses, being faced with lack of informational awareness and with 'information gaps' (Crain, 2018, p. 91), need to infer data and predict behaviours using analytics and modelling to fill that gap. If these sources are 'data poor', the inferences and algorithmic identities (poorly) built on them will be inevitably inaccurate, affecting further the automated decision-making processes.
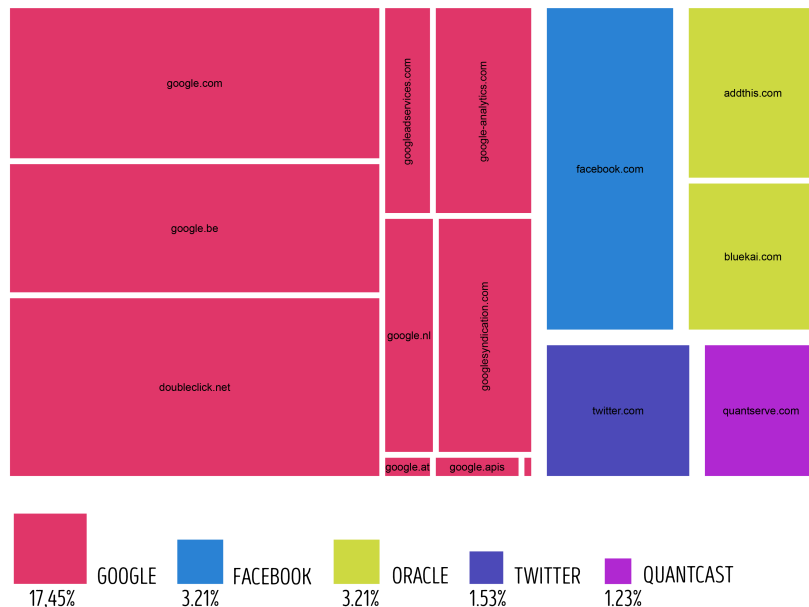


*Figure 2. The most prevalent trackers per company in the research dataset as captured by TrackingObserver, data triangulated with WhoTracksMe and Better.fyi[11]*

---

[11] For the analysis and the visualisation, we focused only on the most prevalent trackers by company, in order to detect the most dominant ones collecting behavioural data about the research subject. That is the reason why the percentages don't sum up to 100% and the long tail is not fully visualised.

*Privacy Score*[12] provided us with different insights. Aiming to investigate which are the websites that capture the user's habits, what kind of trackers are present and for what purposes, we scanned the top 10 most visited websites by the research subject. As detected with the *TrackerObserver*, if we look at the presence of company trackers in the sampled websites, here we also encounter a well-developed network, dominated by *Google* and distantly followed by *Amazon, Oracle, Facebook, Conde Nast and Quantcast* (Figure 3). The analysis further shows that most of the trackers set by third parties are via cookies (73,41%) and for the purpose of advertising (83,23%) (Figure 4). Cross-referencing data collected via *Privacy Score* with data from *Better.fyi* and *WhoTracksMe* enabled us to detect the purpose for tracking and the tracking type detected (Figure 4). This additionally confirms that most of the surveillance done online is for the purpose of accumulating data for online behavioural advertising, referring to personalised and targeted advertising based on prediction of interests and affinities profiling.
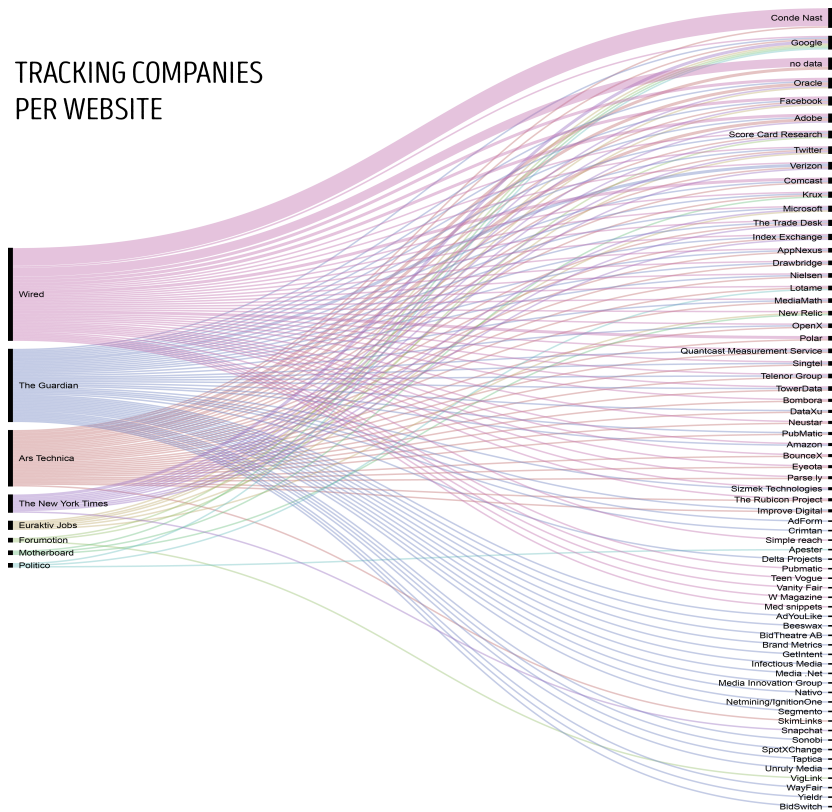


*Figure 3. Frequency of third-party trackers per website. Data source: Privacy Score*

---

[12] Data collected with Privacy Score was done one-time only, as the presence of trackers is tied with the website, not the research subject. The data collection was done on February 21, 2019 and it reflects the state of the particular website at that particular data. Data collected showed that the number of third party embeds (third parties that provide services to the first party) is 575 for only ten websites, set by 328 unique companies, and the number of third-party calls is 172.

*Figure 4. Categories of trackers per website. Data source: Privacy Score, data triangulated with WhoTracksMe and Better.fyi*

The insights collected from these two tools guided the subsequent research steps. It was expected that *Google* and *Facebook* would be the most prominent companies. However, observing the not-insignificant presence of data brokers such as *Oracle* and *Quantcast* motivated the further investigation about the data these companies hold about the research subject and the algorithmic identity they assigned. Data brokers are important actors since they

> are businesses whose revenue model revolves around aggregating information about individuals from a variety of public and private sources […] who sell access to the collected data to third parties, including advertisers, marketers, and political campaigns. (Venkatadri et al., 2018, p.1).

We investigated their role and the data they have by looking at what is detectable regarding datafication practices by different actors at an interface level. To do so, we experimented with data from less-traditional sources: the privacy policies of the most dominant tracking companies we detected in the previous step, the transparency tools made available by the actors themselves and the regulatory tools — the Data Access Request mechanisms enabled by Article 15 of the GDPR.

We started with the *privacy policies* as investigation tools. We sampled the following platforms — *Google* and *Facebook* — and two data brokers — *Oracle* and *Quantcast* — detected previously. To get better initial

structured overview, we used the machine learning tool, *PriBot*[13], in order to collect data on (1) what kind of data is being collected about the users and (2) the reasons for data collection. Although privacy policies can be information-rich sources, we decided to narrow our analysis to these two aspects only, as they are the most relevant for our research question.

## Description of data collected

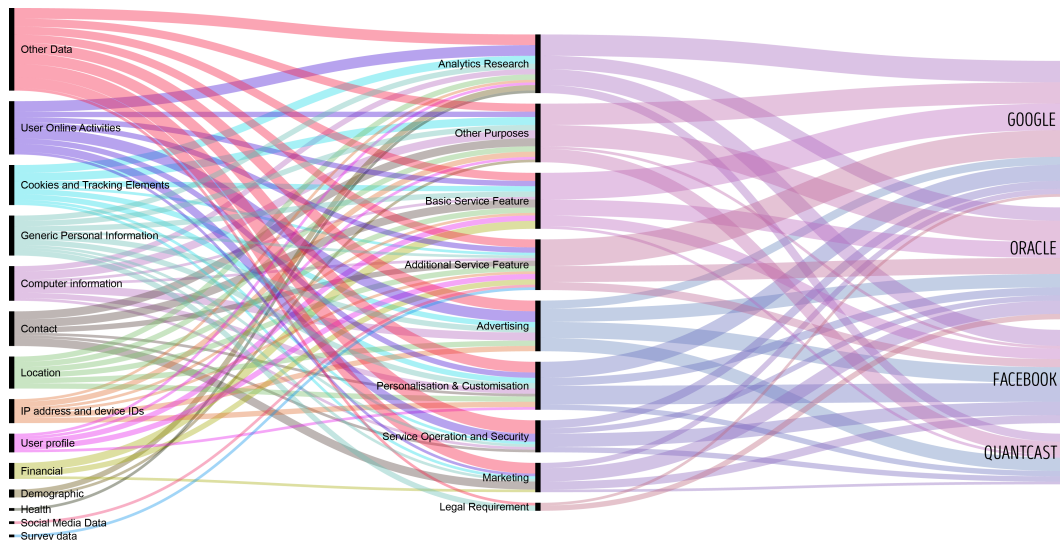| Data | Description | Data | Description |
|---|---|---|---|
| Computer information | The type of operating system (OS) or web browser that the user uses, or similar computer or device information | Location | Geo-location information (e.g. user's current location) regardless of granularity i.e. could be exact location, ZIP code, city-level. |
| Contact | Contact information, such as name, email address, phone number, street address etc. | User profile | The user's profile on the first party website/app, and its contents, e.g. data in user profile, data that user uploaded to website, user comments, user profile preferences, etc. This is common for websites/apps where users can create an account or profile, e.g. on Twitter, YouTube, Facebook, Amazon, etc. |
| Cookies and Tracking Elements | Identifiers locally stored on user's device by the company/organization or third-parties including cookies, beacons, or similar that are commonly used to uniquely identify users, but that are not essential to establish a connection with the user's device or to provide a service | User Online activities | The user's online activities on the first party website/app or other websites/apps, e.g. pages visited, time spent on pages, general user behavior online etc. |
| Demographic | Demographic information, e.g. gender, age, occupation, education, etc. | Social Media Data | User profile and data from a social media website/app or other third party service to which the user gave the first party access, e.g. by connecting with Facebook, Twitter, or other services. Exchanged data may include user profile, photos, comments, friends etc. |
| Financial | Financial information, such as credit/debit card data, other payment information, credit scores, etc. | Survey Data | Any data that is collected through surveys. |
| Health | Health information, such as information about health conditions, prescriptions, medications, as well as health monitoring data, e.g heart rate, step count, activity level etc. | Generic Personal Data | No specific type of information is mentioned, but the policy talks about "personal information" or "personal identifiable information" in general. |
| IP address and device IDs | Permanent (e.g. device IDs, MAC address) or temporary (e.g. IP address) identifiers needed to establish a connection for the current browsing session. | Other data | The type of information is not explicitly stated or unclear (e.g. refers to "information" very generically) |

*Figure 5. Overview of the type of data and specification of data types in the sampled privacy policies (information source: PriBot)*

By analysing and comparing the information we obtained from the *PriBot* tool, a list of all possible data types that could be collected by these actors,

---

[13] *PriBot* is an AI-powered tool for automated analysis of privacy policies https://pribot.org/polisis

the above table was created (Figure 5), listing all the data and their definitions, that were/could be captured for the research subject and further (re)used, and which are of particular interest to the sampled companies. This reveals a dominant 'collect all' approach, where the (legal) principle of data minimization is not respected and a lot of data that is not necessary for establishing a connection or providing a service is captured.

We additionally analysed and cross-referenced the findings for each of the actors (Figure 6), thus being able to discover the relations between the types of data collected by each of the sampled policies, the stated reason for collection and the actors that collect each type. The analysis shows that the most under-defined category — "Other data" is the most frequently captured data, although it was not explicitly stated in any of the policies what kind of data that is, leaving many open doors for misuse and abuse. Looking at the column with particular actors, it is noticeable that apart from *Google* (not unexpected), *Oracle* is actually the actor that closely follows *Google* for potential capture of a number of various data types. Figure 6 shows how "messy" data collection is, and how different types of data can be used for various purposes. We can detect, for example, that *Personalisation & Customisation* is a reason for data collection for all sampled companies, and the following types of data are used for that purpose: user profile, IP address and device IDs, location, contact information, computer information, generic personal information, cookies and tracking elements, user online activities and other data. For *Marketing purposes*, companies use financial data, contact information, generic personal information, cookies and tracking elements, user online activities and other data.



TYPES OF INFO COLLECTED AND REASON FOR COLLECTION PER COMPANY

*Figure 6. Diagram of data collected and stated reasons for collection across companies (information source: PriBot)*

What we call *transparency tools* are designed by the platforms with a specific purpose in mind: to increase transparency and accountability towards users and regulators (Facebook Newsroom, 2020; Google Blog, 2018). However, here we are repurposing them as *objects for study* in order to investigate datafication practices and sources. For this particular case, we looked into *Google* and *Facebook's* data explanation and ads explanation mechanisms.

*Google's* Ad preference page [14], for example, shows the inferred interests about each particular user, briefly elaborating on the logic and process behind it. This allows us to investigate where the (behavioural) data originates from. Having this information, we can see how data is captured and transferred and thus get insights into the datafication and data sharing network. Following and recording the data a few times a week over a period of two and a half months[15], during which we collected 183 distinct interests assigned to the research subject, our research showed that *Google* estimates the interests based on using and/or combining data from: 1. activity on *Google* services/products; 2. activity on *Google* combined with activity on other websites and apps; 3. activity on non-*Google* (outside *Google*) services and 4. Visiting an advertiser's website/app.[16] This also gives insights into the structuring of information and the degree of (non)disclosure by the platform itself, impacting the degree and scope of possible research insights. However, as these systems are highly volatile, at the time of writing this article and checking explanations again, it was noticed that *Google* added one more insight source — "similarity to other users". As an example, for the categorisation "Homeownership Status" *Google* categorises the research subject as "Renter" based on "Google estimates this demographic because your signed in activity on Google services (such as Search or YouTube) is similar to people who've told Google they're in this category". Additionally, three months before, the research subject was categorised as "Homeowner", based on the same sources (see Figure 7).

---

[14] It can be accessed at the following link: https://adssettings.google.com/authenticated

[15] This data was collected in the period of March 2, 2019 to May 17, 2019.

[16] The explanations provided by Google for each of the sources are the following: 1. Google services/products - "Google estimates this interest, based on your activity on Google services (such as Search or YouTube) while you were signed in"; 2. Google and other providers - "Google estimates this interest, based on your signed-in activity on Google services (such as Search or YouTube), as well as on your signed-in activity on non-Google websites and apps"; 3. non-Google (outside Google) - "Google estimates this interest, based on your activity on non-Google websites and apps while you were signed in"; and 4. Visited advertiser - "This advertiser shows you ads based on: Your visit to the advertiser's website/app".
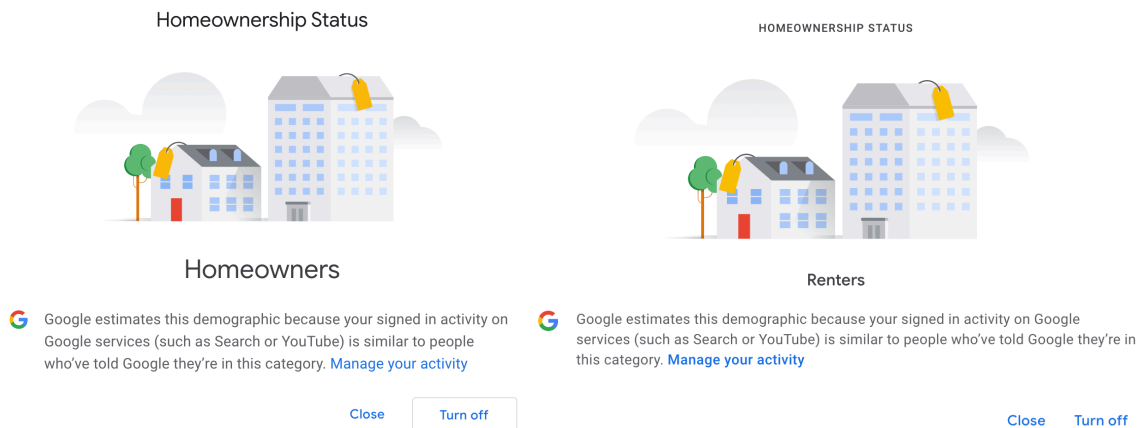
*Figure 7. Screenshots from the research subject's Google's Ad Settings page. The one of the left dates from April 27, 2020; the one on the right is from July 21, 2020*

*Facebook* offers more transparency mechanisms, of which we used the *data explanations*[17] and *ad explanations*[18]. We used these two tools to collect information on the sources of data, the type of data (whether or not personal data) and the actors in the datafication network, as well as — equally important — the mechanisms and sharing practices between the actors in the network.

The insights provided show that *Facebook* datafies users both on- and off-platform, of which the latter one is the prevalent one. Using additional sources of insights about the workings of the platform's tracking system, such as guidelines offered to advertisers by *Facebook* itself, shows that this is data originating from the websites integrating the *Facebook Pixel* tracking technology, and is handed to the platform by clients (websites/app) that integrate it. Clients uploading a contact list to *Facebook* is another source of data feeding the platform. These two sources (*Pixel* and *List*) contain personal data and they constitute 68.25% of the off-platform data ending up at *Facebook*. The only data originating from on-platform behaviour is the data gathered by tracking the ads shown on *Facebook's* Newsfeed that were clicked. Recording the data available via the "advertisers who use contact list added to *Facebook*" tab, shows that very high percent (75%) of companies

---

[17] Data explanations provide the user with a list of attributes Facebook has inferred about them, how they were inferred and what information is used to target them with advertisements (see Andreou et al., 2018 for more detailed explanation of the mechanisms). The data explanations are accessible via an Ad Preferences Page (https://www.facebook.com/ads/settings ) and they provide information structured in the following way: Your interests, Advertisers and Businesses, Your information and Ad Settings.

[18] Ad explanations provide the user with an information/explanation why a particular ad was served. They are accessible via the "Why am I seeing this?" button above every ad served on the user's Newsfeed.

listed collected personal data from other sources, not the user itself, without the user's explicit consent or information about the source provided. This potentially points to the well-developed network of actors in the (personal) data sharing network.

Repurposing the *Ad explanation* tool by *Facebook*, particularly the "Why am I seeing this ad" option, we were able to collect information on the data sources used for personalised ad targeting.[19] We did this on both levels (interface and software), using both observation for recording the data from the interface, and the automated tool *AdAnalyst*, to collect data at a software level. Following an analysis of the explanations provided, we were able to uncover the relations between the sources of data used, the types of data used, the analytical processes at play and the particular reasons for personalised ad targeting, shown in the figure below (Figure 8). For example, if the targeting is based on a *particular interest*, behavioural data will be used to make that inference. This data could be originating either from *Facebook* (by tracking the activity of the user), and advertisers and/or data brokers, using inferential and prediction analytics. The latter analytics methods are used to infer user preferences, attributes and opinions and predict behaviour (Wachter and Mittelstadt, 2018, p.4). Reading, structuring and coding the information collected and recorded, provided us with additional insights: apart from insights into the processes behind the ad-targeting analytics and the inputs/outputs relations, it also revealed that the sources of data could originate both on- and off-platform, they can be volunteered (by the user), obtained (via partners, data brokers and advertisers) or captured by *Facebook*. Different types of data are taken as signals for affinities/interests. This ranges from location and age, to languages spoken, activities and social neighbourhood, or tracking the social network of/for relations between individuals/users and taking this as a data signal for further affinity profiling and commodification for ads targeting. This 'data inference process' (Andreou et al., 2018, p.3) is important because it allows the advertising platform to infer users' preferences and attributes, later used for affinity profiling and building algorithmic identity, further used as a basis for commodification (targeted advertising).

---

[19] Such as: liked advertising page, visited advertiser's website or app, friends liked a page, age/gender/location, activity on Facebook's family of apps & service, particular interest etc.

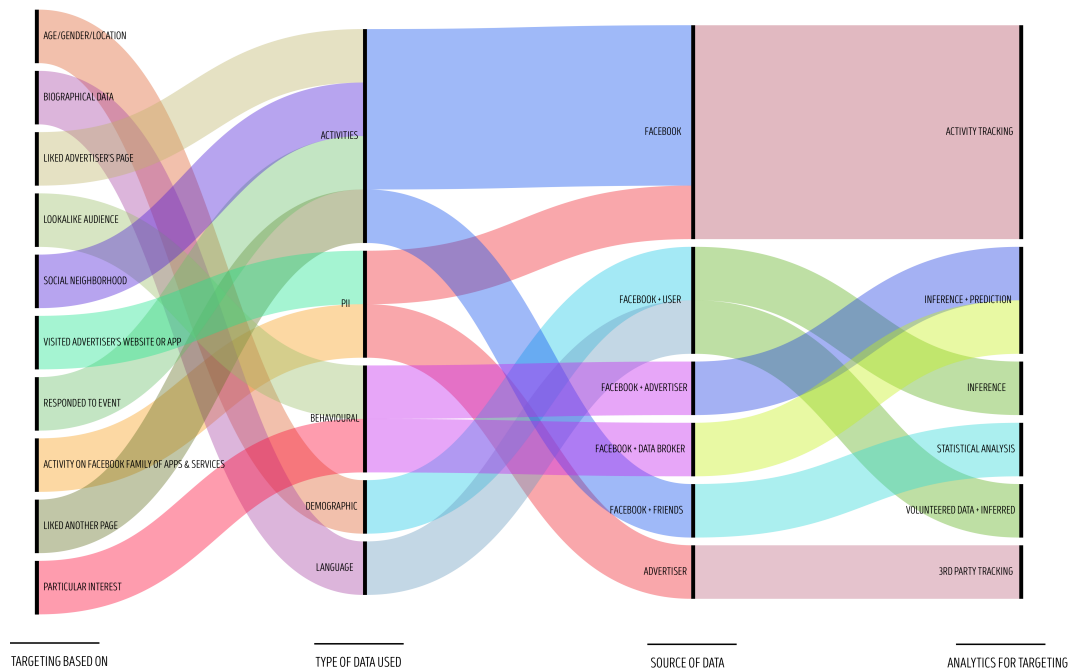FACEBOOK TARGETING BY TYPE, SOURCE AND ANALYTICS OF DATA



*Figure 8. Alluvial diagram of sources of data and inferences for Facebook*

The last strategy we used for uncovering and investigating the data sources, actors and mechanisms for inferential analytics, prediction and building algorithmic identity was repurposing the *Data Subject Access Rights* mechanisms as an object for study. Article 15 of the GDPR, in force since May 2018, enables data subjects to request and obtain access to any personal data being held and processed by a data controller. Executed in correct manner, it should give information on the purposes of the processing, the categories of personal data concerned, the recipients or categories of recipients to whom the personal data have been or will be disclosed, and if automated decision-making (including profiling) is present. The latter entails providing meaningful information about the logic involved, as well as the significance and the envisaged consequences of such processing for the data subject (European Commission, 2016). Repurposed for academic research, Data Subject Access Requests (DSAR) would give information on the sources of data (categories of personal data concerned), the network of actors with access to the data and the algorithmic identity/assigned affinities by the controller.

Six DSARs were filed, of which only one response (by *Oracle)* was entirely suitable for analysis.[20] The data obtained from *Quantcast*, although

---

[20] Requests were sent to *Bumble*, *Oracle*, *Criteo*, *Quantcast*, *Facebook* and *Acxiom*. Only *Oracle* provided data that can be used for the purposes of the research. The file obtained by *Quantcast* was "unreadable" in terms that it contained only a few unique rows, duplicated tens of thousands of times (96,659 data entries in total). *Criteo* was asking for additional identification checks, and because of time constraints it was decided not to follow through.

incomplete, enabled for some crucial observations. The first observation pertains to the well-established and wide network of data sharing and the exchange system between data brokers. *Oracle* relies on six other data brokers to collect data and infer affinities and interests (these data brokers are *Eyeota, OnAudience, Lotame, Bombora, AuDigent, Affinity Answers*). This complicates the quest of tracking where data originates and where its' final destination is, making it difficult to later contest or rectify the data in question. The second observation concerns the risk of inaccurate inferences: if one data broker makes inaccurate inference, this information is further shared across the ecosystem. Closely inspecting the data provided by *Oracle*, it could be observed that some of the inaccurate data *Oracle* holds originates from *Eyota*, that obtained them from *Bombora*. The reliance on other partners and data brokers is also indicated in the data obtained from *Quantcast*, in their "Audience Grid" data file, which points to a largely adopted practice. This might have serious consequences for the data subject resulting in not just their erroneous profiling, but also (potentially) in access to services and opportunities.

The "unsuccessful" DSARs also demonstrate that the access to personal data held by online platforms is more often than not a complex and uncertain process. Because of the different interpretations of the DSAR procedure and the GDPR in general by companies, there are apparently substantial differences about what data is considered personal and thus eligible to be provided by the data controllers.[21] Sometimes the data controllers have long and extensive procedures (like *Criteo*) or they try to bypass meaningful information by directing users towards other available data (*Facebook*). Even when successful, the data obtained might not be readable (as in the *Quantcast* case), the file might be incomplete, and the logic behind the presented and provides data and information might not be available or accessible for the user.

## 5.2 Investigating algorithmic identity

Next we investigated the workings of the algorithmic systems of a web platform (*Google*), social media (*Facebook*) and one data broker (*Oracle*). We took the inferences as proxies, or represents, for investigating the assigned

---

*Facebook* provided data, but with no additional meaningful information, and the data corresponds with the one provided on their platform via the "Download your information" tool. *Acxiom* provided an answer stating that no data is collected from individuals residing in Belgium.

[21] In an attempt to obtain data from the dating app *Bumble*, the platform representatives stated that they can only provide a registration date, IP addresses and profile photos (source: personal correspondence).

algorithmic identity. We decided for sampling these two platforms and the data broker based on the results from the datafication phase of the research, where most trackers were originating from these three actors (and as such have most data on the research subject), and on their affordances for research.
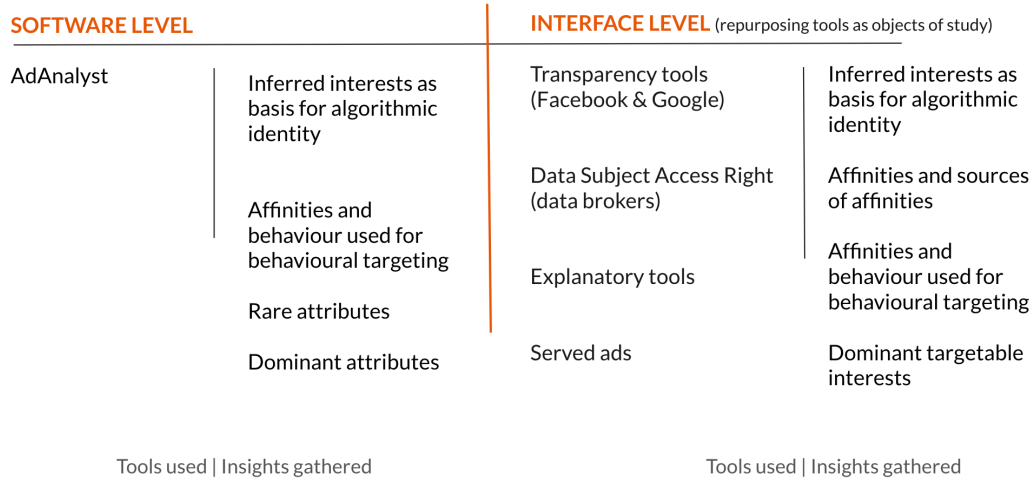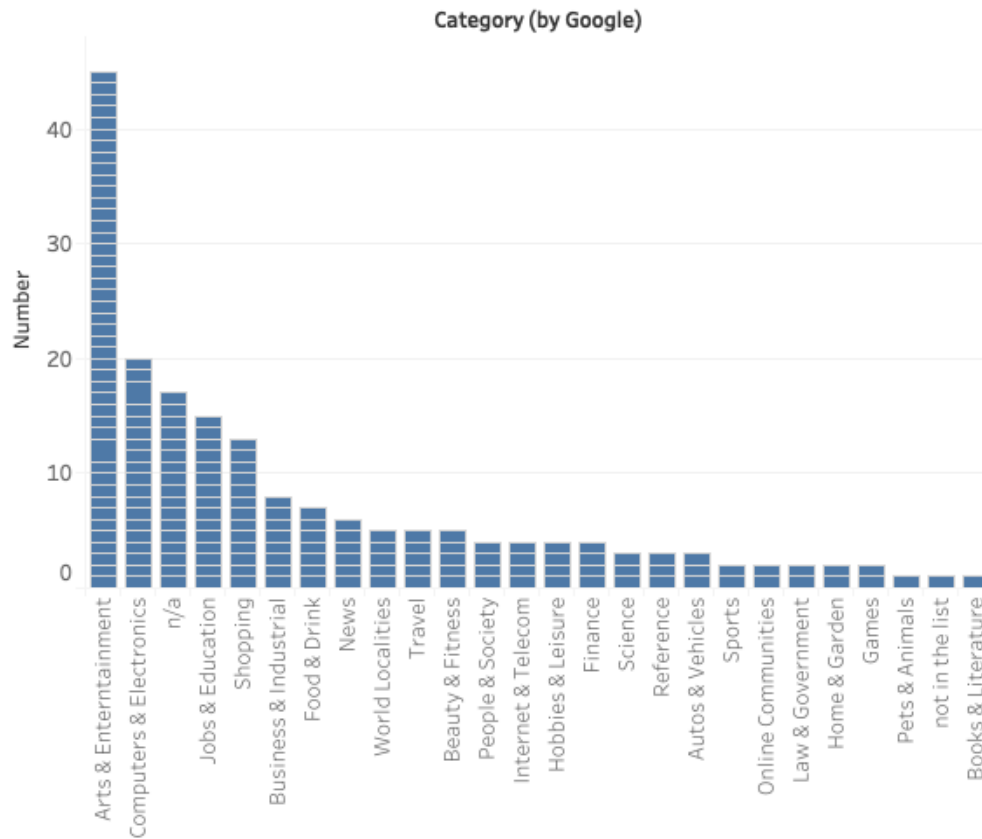
| SOFTWARE LEVEL | | INTERFACE LEVEL (repurposing tools as objects of study) | |
|---|---|---|---|
| AdAnalyst | Inferred interests as basis for algorithmic identity | Transparency tools (Facebook & Google) | Inferred interests as basis for algorithmic identity |
| | Affinities and behaviour used for behavioural targeting | Data Subject Access Right (data brokers) | Affinities and sources of affinities |
| | Rare attributes | Explanatory tools | Affinities and behaviour used for behavioural targeting |
| | Dominant attributes | Served ads | Dominant targetable interests |
| | Tools used \| Insights gathered | | Tools used \| Insights gathered |

*Figure 9. Overview of the tools used for algorithmic identity and insights gathered*

The three different data controllers are investigated in order to assess the assigned algorithmic identity and test the possibility for research using the inferred interests as proxies. As Figure 9 shows, we employed a variety of methods and tools, at a different level, to analyse various aspects of the inferential analytics at play and their outputs.

*Google's Ad settings* tool[22] was observed in frequent intervals for two and a half months and it was used to record the assigned interests. Based on the 184 observed interests, and triangulated with a list of categories (Brave, 2019) indicating the particular category an interest belongs to, enhanced with a close reading of the categories, we were able to get an overview of the most dominant categories the research subject was categorised in (Figure 10). The daily recording of the interactions and assigned interests show that these are often immediate outputs of simple browsing behaviour, but also that they are unstable and disappearing – thus no historical database of inferred interests is available (for research or personal insights). Some of the interests disappear on a daily basis and some remain longer periods of time, or during the entire period of data collection. This is significant from a point of view of reliability of collected data: researchers must be aware of the instability of the data and the potential inability of collecting what is available. This underlines the

---

[22] The information available by the platform was monitored, collected and recorded in the time period of March 2, 2019 – May 17, 2019, and 183 interests assigned were observed.

dependence on and significance of the information structuring and information visibility, which can be seen as politics both of visibility and knowledge, controlled by the platforms themselves. Andreou et al. (2018) point to the same characteristic of *Facebook's* transparency tools, referring to it as *snapshot/temporal completeness*.



Sum of Number for each Category (by Google). Details are shown for Interest.

*Figure 10. Frequency of categories of interests as assigned to the research subject by Google*

Reading the assigned interests as *text*, we were able to construct an overview of the assigned algorithmic identity by *Google* (Figure 11). The use of the auto-technographic approach, as well as the fact that we are relying on and working with data from a real individual, enables us to test the assumptions made by the algorithms and assess its truthfulness. In our case, the assigned algorithmically constructed identity is in a sharp discrepancy with the research subject's sense of real identity and does not represent their actual life conditions (financial, familial, or employment). Similar are the findings from the data collected from *Oracle*, with the important difference here that online data brokers often lack information on basic demographic data and thus have to infer it via browsing behaviour in order to fill the 'information gap' (Crain, 2018), unlike platforms like *Google* and *Facebook*

that rely on both volunteered data (by users) and have more access to daily behaviour of users. However, we must be aware as researchers that an important aspect of reading and interpreting the data is concealed by the platforms: there is a lack of information on how these attributes are assigned, and what is the inferential analytics process. This potentially affects the comprehensiveness of the data collected by the researcher and consequently — the analysis itself.
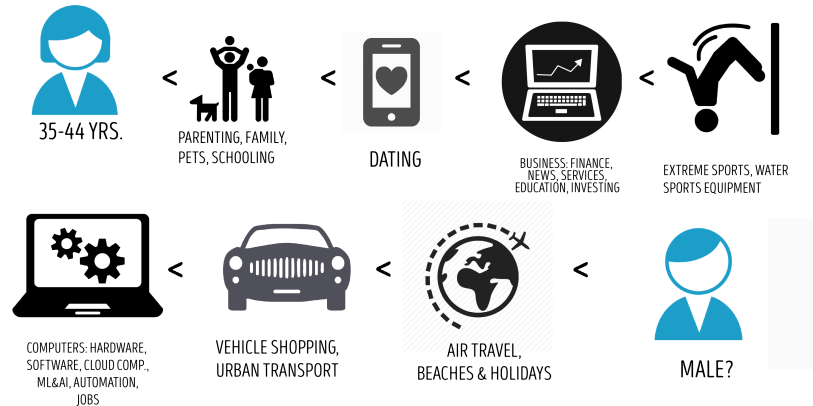
## GOOGLE ASSIGNED ALGORITHMIC IDENTITY

35-44 YRS.

PARENTING, FAMILY, PETS, SCHOOLING

DATING

BUSINESS: FINANCE, NEWS, SERVICES, EDUCATION, INVESTING

EXTREME SPORTS, WATER SPORTS EQUIPMENT

COMPUTERS: HARDWARE, SOFTWARE, CLOUD COMP., ML&AI, AUTOMATION, JOBS

VEHICLE SHOPPING, URBAN TRANSPORT

AIR TRAVEL, BEACHES & HOLIDAYS

MALE?

*Figure 11. A close-reading illustration of an algorithmic identity as assigned by Google*

## ORACLE ASSIGNED ALGORITHMIC IDENTITY

18+/65+ YRS.

EMPLOYED + RETIRED

DECLARED DAD

BUSINESS & IT PROFESSIONAL

SPORTS

INVESTMENT & STOCK EXCHANGES

PET OWNER

GOVERNMENT & POLITICS ENTHUSIAST

GREEN LIVING & ENVIRONMETALISM

*Figure 12. A close-reading illustration of an algorithmic identity as assigned by Oracle*

When it comes to the possibility to investigate algorithmic identity as assigned by *Facebook*, by using the very affordances of the platform itself, we were able to draw an overview of the assigned general affinity towards certain categories, via few available "points". We used as data source the

*data explanations* (revealing the reason for assigning the interests[23]) and the *ad explanation* feature, both at an interface level (via observation and recording of data) and at software level (*AdAnalyst* tool).

The data collected via the data explanation feature, gives not only insights into the dominant assigned interests by category (Figure 13), but also points to the very specific categorisation practices *Facebook* uses for profiling and targeting. Closely reading the list of interests, it becomes visible that *Facebook* is constructing very narrow categories (e.g. *headphones; old style and new style dates; Conversion (gridiron football); Right-to-work law*; particular movies/songs etc.) that might enable a very specific targeting, and also, that many of them simply do not make sense (e.g. *non-resident Indian and person of Indian origin; hydrogen*) or can be regarded as potentially sensitive information (*Gay Pride, LGBT community*).



INFERRED FACEBOOK INTERESTS BY CATEGORY

*Figure 13. Dominance of interests assigned by Facebook, per category*

Although the matching process of a user being served a particular ad is complex due to the fact that the outcome doesn't depend only on the advertising platform and its matching algorithm, but also on the very event-specific factors[24], the explanatory tool "Why am I seeing this ad?", when

---

[23] At the time of the data collection and analysis of the *Facebook* data (11/27/2018 - 04/16/2019), *Facebook* was providing three very generic explanations why an interest/affinity was inferred, but no further information: "you clicked an ad related with the interest" (64.37%), "you liked a page related with the interest" (31.80%) or "installed an app" (0.5% of the entries). Additionally, it added "liked their page or post" (3.30%) - data recorded in May 2019 showed that *Facebook* made changes to their "assigned interests" explanation, adding one more "reason" to the previous three.

[24] Such as the competing advertisers at the particular moment when an ad is about to be served, their specific requests/objectives set by the advertisers and the characteristics of the available users on the platform, in a particular moment of time (Andreou et al., 2018, p. 3).

repurposed as an object for study, can provide significant information regarding the particular behaviour, activities and interests of the research subject used for automated behavioural targeting. Combining the insights collected manually via the interface with the data collected automatically via *AdAnalyst* at a software level, provided significant insights. The first finding is related with the type of data that algorithmic systems consider as an important particular aspect of the research subject's algorithmic identity to be later taken into account for personalised behavioural ad targeting.[25] The second one relates to the affordances of the different research methods and tools, and the different insights, depth and scope of insight that they enable. *AdAnalyst* offers different insights as it has access to more parameters at a software level, not accessible via the interface. Such is the distinction between the general ad explanation served to the research subject (as a user) and what is indicated as a reason the particular user to be targeted. Additionally, insights can be obtained about the targeting parameters set by the advertisers. As Figure 14 shows, what the research subject has been targeted based on (e.g. bicycle as interest), might be just one of the campaign targets set by the advertiser. These can sometimes be different, and in that sense *AdAnalyst* provides more in-depth insights than available if looking only at an interface level.



*Figure 14. Screenshot of AdAnalyst's interface*

The screenshot above is interesting for analysis because, via the section "The advertisers targeted other users with", it provides valuable insights into the parameters *Facebook* uses for targeting. We can observe that apart from the well-known Lookalike audience, Personally Identifiable Information (PII)[26], Social Neighbourhood and similar, it also targets users

---

[25] For example, the research subject has liked a page, has or was at a particular location, belongs within a particular age group, etc.

[26] Personally Identifiable Information (PII) is considered any data that can be used to identify a specific individual, such as name, email, phone number, IP address, location address, online identifier, biometric records and similar. For more detailed definition, see GDPR Art. 4 (1).

based on data from data brokers, based on behaviours (e.g. expats in France), operating system and version (based on where *Facebook* was accessed from) and biographical data (Master's degree).

Another avenue to investigate and assess an assigned algorithmic identity is to repurpose the particular ads served to the user, more specifically the textual part of each of the ads. As the purpose of the ads is to nudge users to take particular action, ads are served targeting specific interests of particular users, with the aim to steer actions or behaviour. In that sense, ads could uncover the assigned affinities and, at an aggregate level, the algorithmic identity. Thus, a semantic analysis of 1,553 served ads, collected both manually (interface level) and using *AdAnalyst* (software level), was done. Only unique ads were taken into consideration. The tool *CorText* (Munk, 2019) was used to detect the semantic clusters forming from the corpus of served ads. The frequency of the semantic co-occurrence can be read as a signal of attributes the user is more targetable for, or most prone to take actions for. It can also be seen as enabling an insight into how a particular user is seen by the algorithms, given that the most dominant reasons for targeting are being part of a lookalike audiences and because of specific user interests.
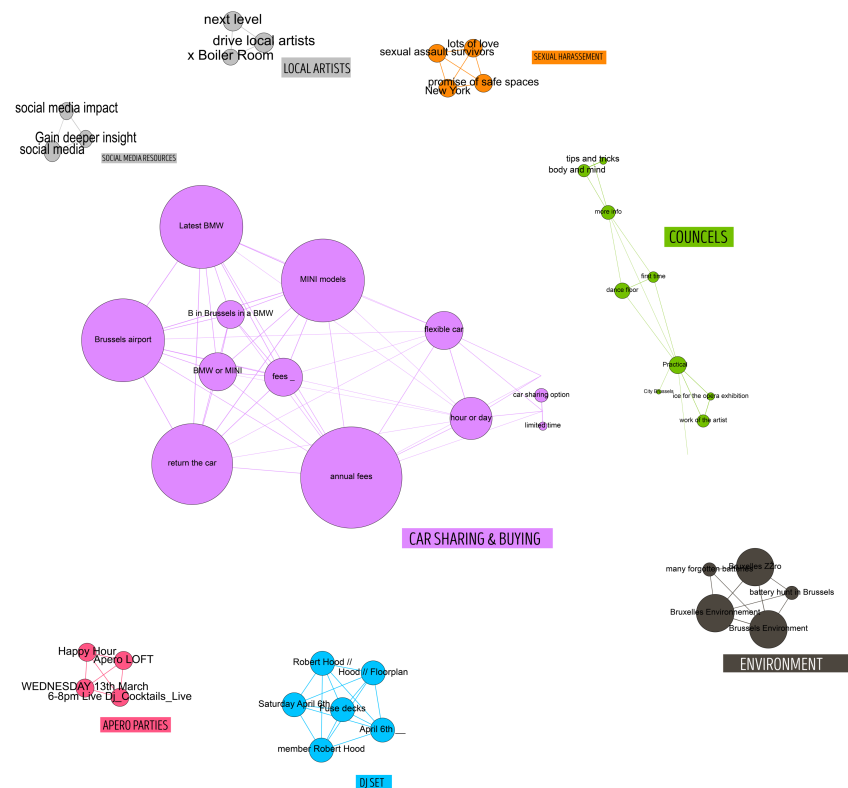
## SEMANTIC ANALYSIS OF THE FACEBOOK ADS SERVED



*Figure 15. Network mapping of semantic clusters from served ads on Facebook, using CorText*

111

# 6    CONCLUSION: ONE APPROACH TO GUIDE THEM ALL?

As Marres and Gerlitz (2015) observe, social media platforms 'do not present us with raw data, but rather with specially formatted information' (p.22). The formatting of this data, both at an interface and software (API) level, then inevitably influences the methodological implications for research. By "standardising" the presentation of data and the way it is made visible, the platforms are guiding the researchers through what is available to be seen and investigated. The perspective, methods and insights are limited by the affordances of each and every platform, their algorithmic and API system(s). Keeping this in mind is important for discussing the scope and depth of available information when employing the set of research methods and tools in this empirical research. Marres and Gerlitz (2015) call this 'methodological bias' (p. 22) and rightfully ask the question if "is it really the researcher that here 'decides' to use this method, or is this decision rather informed by the object of study with its associated tools and metrics?" (ibid.). If not limited in the right sense of the word, then we, as researchers, are nudged, steered towards the particular configuration of analytic practices via the platforms', APIs' and software's own 'sampling techniques, options for analysis and modes of visualization' (p. 31).

Another potentially problematic aspect of relying both on APIs and data and being denied access to them and adopting a method for data collection based on observation, is the constant change of what platforms make available. This highlights the constant revision and change of their *politics of visibility* and *politics of knowledge* implemented via the changes at an interface and software level. Barrett and Kreiss (2019) call this *platform transience* – a concept they use to describe the sudden changes platforms make in their policies, procedures and/or affordances, which impacts the ability for critical research, as it makes them continuously changeable and ephemeral in significant ways. Right after the end of the data collection phase of this research, *Facebook* changed its data and ad explanation structures, now offering more information at the disposal of users (Facebook Newsroom, 2019). This is not just problematic in the sense that it makes data collected at different time-periods potentially incomparable, but it also makes the study of algorithmic systems almost a study of 'historical objects' (Bucher, 2012). We as researchers will be always bounded by what platforms decide to make available, either via the interface or the API. With platforms closing their APIs and giving data access only to the "chosen" few (for example, *Facebook's Social Science One*[27]), move described by Bruns (2019) as 'corporate data philanthropy', the data access gap will

---

[27] More about the *Facebook*-Academic partnership can be read at the following link: https://socialscience.one/. Last accessed July 31, 2020.

be only widening. Hence our ability to study technology, society, and the intersection of the two, will narrow down and become potentially very limited.

Considering this, and considering the increasing limitations of how research can be done and what can be obtained as valuable knowledge, as a result of immanent methodological bias, API restrictions and impenetrability of black boxes, we are faced with the question of how successful and valid the research we conducted was. At an interface level, the methodological design imposed some limitations in a particular manner. This is once again related with how the platforms organise their information: how *Facebook* and *Google's* ad settings are organised, how much they reveal, how the data obtained via Data Subject Access Request is organised, how readable it is and finally, what is made available through these "interfaces" and what is concealed, left out or not provided. Another related aspect concerns the nature of observation as a research method and of the auto-technographic approach. As outlined by Weltevrede (2016), this is always a *me-centric* view, highly individualised and personalised (p. 107), as is our experience and content provided on these platforms. Additionally, we need to be aware of the complexities arising when one would like to translate this very same methodological design and setting on a sample comprising more than one research subject. That would require undertaking additional and modified steps, setting up the research environment and testing the possibilities to obtain valuable and valid data, considering all the complexities of browsing histories, browsing habits and patterns, that particular research subjects could exhibit.

These exact same limitations can be seen as an advantage, as they enable 'real user-algorithmic agent interactions' (Bodo et al., 2018, p. 143). Being able to observe these enriches the quality of the insights, but more importantly, it allows to see the wider 'socio-technological assemblage' (ibid.) and the networks between different actors. And while it might not provide a picture of the totality of the system, it does provide a valuable, although partial, reconstruction of the complexity of these algorithmic assemblages.

By using the affordances of the different methods, at a different level of visibility (interface and software) for analytical inquiry, and combining these findings, new and more in-depth insights were made possible. This is reinforced with the action of repurposing objects of/for study — such as the data explanations, ad explanations, data subjects access request and similar — as a strategy to overcome the limitations, uncover and make visible what was previously not *revealable*. While having to adjust to the affordances and thus limitations of methods and tools, this research and methodological strategy offered ways to be innovative, to — by learning what is possible — look for new avenues, new perspectives, new sources of data and thus

insights for digital social research. In that regard, the methodological design of this research is successful, as it provides access to new insights and enables for a more in-depth inquiry into the processes of algorithmic construction of identity, data extraction and inferential analytics, and the ecosystem of actors and networks around these surveillance practices. At a software level, automated tools enabled for a more in-depth knowledge and helped better investigate aspects hidden from the interface and the eye. However, the approach has its limitations, emanating from the nature of platforms' APIs, which are also limited in scope and applicability by their very affordances. They have their own "politics of visibility", limiting what can be seen and uncovered. At an interface level, the daily, detailed observation and recording of the workings and outputs of the system enable for more granular insights and observations of the subtle changes in and by algorithmic systems.

With our research we tried to manoeuvre around the restrictions for research imposed by APIs and black boxes and find ways to investigate opaque algorithmic systems. Following Paßmann and Boersma's (2017) suggestion for pursuing practical transparency, complemented by what they call formalized transparency, we made use of sources external to the algorithms, their APIs and black boxes as a way to detect and make known the unknowns. While APIs are important research entry point, they are not the only one. We experimented with different approaches to circumvent the limitations for research imposed by platforms' gatekeeping practices. In doing so, we got close to what can be called 'digital fieldwork' (Venturini and Rogers, 2019): exploring, experimenting with, testing and employing various new approaches, sources, ways to collect data and capture the interactions between the algorithms and users, mediated via interfaces and APIs. With that, we proposed (just) one of the possible avenues for overcoming data access gaps and algorithmic opacity in doing digital social research. While the question of *if* and *how* platforms should provide access to data for researchers is not a focus of this paper, it remains an important one. We are on the opinion that while it is necessary, thorough digital social research should use and rely on other methods, techniques and data access points in combination with API data. We see this as the only approach that will provide comprehensive view of the socio-technological assemblages, their outputs and impact.

## FUNDING STATEMENT AND ACKNOWLEDGMENTS

## REFERENCES

*AdAnalyst: Bringing transparency to Facebook Ads*. (2019, June 12). https://adanalyst.mpi-sws.org/#about-transpad

Ananny, M., & Crawford, K. (2018). Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability. *New Media & Society*, *20*(3), 973–989. https://doi.org/10.1177/1461444816676645

Andrejevic, M., & Gates, K. (2014). Big Data Surveillance: Introduction. *Surveillance & Society*, *12*(2), 185–196. https://doi.org/10.24908/ss.v12i2.5242

Andreou, A., Venkatadri, G., Goga, O., Gummadi, K. P., Loiseau, P., & Mislove, A. (2018). *Investigating ad transparency mechanisms in social media: A case study of Facebook's explanations*. http://www.eurecom.fr/publication/5414

Angwin, J., Mattu, S., Paris, Jr. Terry. (2016, December 27). *Facebook Doesn't Tell Users Everything It Really Knows About Them* [Text/html]. ProPublica. https://www.propublica.org/article/facebook-doesnt-tell-users-everything-it-really-knows-about-them

Ariana Tobin, J. B. M. (2018, September 18). *Facebook Is Letting Job Advertisers Target Only Men* [Text/html]. ProPublica. https://www.propublica.org/article/facebook-is-letting-job-advertisers-target-only-men

Barrett, B., & Kreiss, D. (2019). Platform transience: Changes in Facebook's policies, procedures, and affordances in global electoral politics. *Internet Policy Review*, *8*(4). https://policyreview.info/articles/analysis/platform-transience-changes-facebooks-policies-procedures-and-affordances-global

Bashyakarla, V. (2018a, May 18). *Psychometric Profiling: Persuasion by Personality in Elections*. Our Data Our Selves. https://ourdataourselves.tacticaltech.org/posts/psychometric-profiling/

Bashyakarla, V. (2018b, September 11). *Geotargeting: The Political Value of Your Location*. Our Data Our Selves. https://ourdataourselves.tacticaltech.org/posts/geotargeting/

BBC News. (2018, October 23). *Mobile app data sharing "out of control."* https://www.bbc.com/news/technology-45952466

Beckett, L. (2014, June 13). *Everything We Know About What Data Brokers Know About You* [Text/html]. ProPublica.

https://www.propublica.org/article/everything-we-know-about-what-data-brokers-know-about-you

Better. (2019, July 18). *Trackers Collection*. https://better.fyi/trackers/

Binns, R., Lyngs, U., Van Kleek, M., Zhao, J., Libert, T., & Shadbolt, N. (2018). Third Party Tracking in the Mobile Ecosystem. *Proceedings of the 10th ACM Conference on Web Science*, 23–31. https://doi.org/10.1145/3201064.3201089

Binns, R., Zhao, J., Kleek, M. V., & Shadbolt, N. (2018). Measuring Third-party Tracker Power Across Web and Mobile. *ACM Trans. Internet Technol.*, *18*(4), 52:1–52:22. https://doi.org/10.1145/3176246

Bodo, B., Helberger, N., Irion, K., Borgesius, K. Z., Moller, J., Velde, B. van de, Bol, N., Es, B. van, & Vreese, C. de. (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.

Boerman, S. C., Kruikemeier, S., & Borgesius, F. J. Z. (2017). Online Behavioral Advertising: A Literature Review and Research Agenda. *Journal of Advertising*, *46*(3), 363–376. https://doi.org/10.1080/00913367.2017.1339368

Boyd, D., & Crawford, K. (2012). Critical questions for big data: Provocations for a cultural, technological, and scholarly phenomenon. *Information, Communication & Society*, *15*, 662–679.

Bozdag, E. (2013). Bias in Algorithmic Filtering and Personalization. *Ethics and Inf. Technol.*, *15*(3), 209–227. https://doi.org/10.1007/s10676-013-9321-6

Brave. (n.d.). *Marked up copy of Google's RTB "Publisher Verticals."* https://brave.com/wp-content/uploads/2019/01/Google-publisher-verticals-marked-up.pdf

Bruns, A. (2019). After the 'APIcalypse': social media platforms and their fight against critical scholarly research, Information, Communication & Society, 22:11, 1544-1566, DOI: 10.1080/1369118X.2019.1637447

Brusseau, J. (2019). Ethics of identity in the time of big data. *First Monday*, *24*(5). https://doi.org/10.5210/fm.v24i5.9624

Buchanan, R. (2001). Design research and the new learning. Design Issues. 17(4), 3-23.

Bucher, T. (2012). Want to be on the top? Algorithmic power and the threat of invisibility on Facebook. *New Media & Society*, *14*(7), 1164–1180. https://doi.org/10.1177/1461444812440159

Bucher, T. (2016). Neither Black Nor Box: Ways of Knowing Algorithms. In S. Kubitschko & A. Kaun (Eds.), *Innovative Methods in Media and Communication Research* (pp. 81–98). Springer International Publishing. https://doi.org/10.1007/978-3-319-40700-5_5

Bucher, T. (2017). The algorithmic imaginary: Exploring the ordinary affects of Facebook algorithms. *Information, Communication & Society*, *20*(1), 30–44. https://doi.org/10.1080/1369118X.2016.1154086

Bucher, T. (2018). *If...Then: Algorithmic Power and Politics*. Oxford University Press.

Bucher, T. (n.d.). *Programmed sociality: A software studies perspective on social networking sites*. 221.

Cabañas, J. G., Cuevas, Á., & Cuevas, R. (2018). Facebook Use of Sensitive Data for Advertising in Europe. *ArXiv:1802.05030 [Cs]*. http://arxiv.org/abs/1802.05030

Cheney-Lippold, J. (2011). A New Algorithmic Identity: Soft Biopolitics and the Modulation of Control. *Theory, Culture & Society*, *28*(6), 164–181. https://doi.org/10.1177/0263276411424420

Crain, M. (2018). The limits of transparency: Data brokers and commodification. *New Media & Society*, *20*(1), 88–104. https://doi.org/10.1177/1461444816657096

Crawford, K., Lingel, J., & Karppi, T. (2015). Our metrics, ourselves: A hundred years of self-tracking from the weight scale to the wrist wearable device. *European Journal of Cultural Studies*, *18*(4–5), 479–496. https://doi.org/10.1177/1367549415584857

Crawford, K. (2013, April 1). The Hidden Biases in Big Data. *Harvard Business Review*. https://hbr.org/2013/04/the-hidden-biases-in-big-data

Dance, G. J. X., LaForgia, M., & Confessore, N. (2018, December 18). As Facebook Raised a Privacy Wall, It Carved an Opening for Tech Giants. *The New York Times*. https://www.nytimes.com/2018/12/18/technology/facebook-privacy.html

de Montjoye, Y.-A., Hidalgo, C. A., Verleysen, M., & Blondel, V. D. (2013). Unique in the Crowd: The privacy bounds of human mobility. *Nature Research Journal*, *3*, 1376. https://doi.org/10.1038/srep01376

Desjardins, J. (2018, November 19). *Here's What the Big Tech Companies Know About You*. Visual Capitalist. https://www.visualcapitalist.com/heres-what-the-big-tech-companies-know-about-you/

Dias, T., & Natusch, I. (2017, November 29). They are stalking you to calculate your score. *Chupadados*. https://chupadados.codingrights.org/en/they-are-stalking-you-to-calculate-your-score/

Dieter, M., & Tkacz, N. (2020). The Patterning of Finance/Security: A Designerly Walkthrough of Challenger Banking Apps. *Computational Culture*, *7*. http://computationalculture.net/the-patterning-of-finance-security/

Digital Methods Initiative. (2019, July 27). *The research browser < Dmi < Foswiki*. https://wiki.digitalmethods.net/Dmi/FirefoxToolBar#The_research_browser

Dourish, P. (2016). Algorithms and their others: Algorithmic culture in context. *Big Data & Society*, *3*(2), 2053951716665128. https://doi.org/10.1177/2053951716665128

Facebook. (2019, May 21). *How does Facebook detailed targeting work?* Facebook Ads Help Center. https://www.facebook.com/business/help/182371508761821

Facebook Business. (2019, June 13). *About Lookalike Audiences for Facebook ads*. Facebook Ads Help Center. https://www.facebook.com/business/help/164749007013531

Facebook Business. (2019, June 24). *Lookalike Audiences*. Facebook Business. https://en-gb.facebook.com/business/learn/facebook-ads-lookalike-audiences

Facebook Newsroom. (2019, July 11). *Understand Why You're Seeing Certain Ads and How You Can Adjust Your Ad Experience*. https://newsroom.fb.com/news/2019/07/understand-why-youre-seeing-ads/

Facebook Newsroom. (2020, March 30). *Updating Our Data Access Tools*. https://about.fb.com/news/2020/03/data-access-tools/

Fix AdTech. (2019, February 20). *New evidence filed in RTB complaint*. Fix AdTech. https://fixad.tech/february2019/

Fritsch, K. (2018). Towards an Emancipatory Understanding of Widespread Datafication. *SSRN Electronic Journal*. https://doi.org/10.2139/ssrn.3122269

Ghostery team. (2017). *Tracking the Trackers: Analysing the global tracking landscape with GhostRank*. https://www.ghostery.com/study/

Gillespie, T. (2012, February 2). *Can an Algorithm be Wrong?* Limn. https://limn.it/articles/can-an-algorithm-be-wrong/

Goga, O. (2019, May 15). *Facebook's "transparency" efforts hide key reasons for showing ads*. The Conversation. http://theconversation.com/facebooks-transparency-efforts-hide-key-reasons-for-showing-ads-115790

Golebiewski, M., & Boyd, D. (2018, May 11). Data Voids: Where Missing Data Can Easily Be Exploited. *Data & Society*. https://datasociety.net/output/data-voids-where-missing-data-can-easily-be-exploited/

Google Blog (2018, June 14) *Greater transparency and control over your Google ad experience*. https://blog.google/technology/ads/greater-transparency-and-control-over-your-google-ad-experience/

Gutwirth, S., & De Hert, P. (2008). Regulating Profiling in a Democratic Constitutional State. In M. Hildebrandt & S. Gutwirth (Eds.), *Profiling the European Citizen: Cross-Disciplinary Perspectives* (pp. 271–302). Springer Netherlands. https://doi.org/10.1007/978-1-4020-6914-7_14

Haggerty, K. D., & Ericson, R. V. (2000). The surveillant assemblage. *The British Journal of Sociology*, *51*(4), 605–622. https://doi.org/10.1080/00071310020015280

Hill, K. (2018, September 26). *Facebook Is Giving Advertisers Access to Your Shadow Contact Information*. Gizmodo. https://gizmodo.com/facebook-is-giving-advertisers-access-to-your-shadow-co-1828476051

Jackson, V., Rosenberg, D., Williams, T. D., Brine, K. R., Poovey, M., Stanley, M., Garvey, E. G., PhD, M. K., Raley, R., Ribes, D., Jackson, S. J., & Bowker, G. C. (2013). *"Raw Data" Is an Oxymoron* (L. Gitelman, Ed.). The MIT Press.

Jarrett, K. (2014). A Database of Intention? In R. König & M. Rasch (Eds.), *Society of the query reader: Reflections on web search*. Inst. of Network Cultures.

Just, N., & Latzer, M. (2017). Governance by algorithms: Reality construction by algorithmic selection on the Internet. *Media, Culture & Society*, *39*(2), 238–258. https://doi.org/10.1177/0163443716643157

Kaltheuner, F. (2018, November 7). *I asked an online tracking company for all of my data and here's what I found*. Privacy International. http://privacyinternational.org/feature/2433/i-asked-online-tracking-company-all-my-data-and-heres-what-i-found

Kien, G. (2008). Technography = Technology + Ethnography: An Introduction. *Qualitative Inquiry*, *14*(7), 1101–1109. https://doi.org/10.1177/1077800408318433

Koebler, J., & Maréchal, D. N. (2018, November 16). Targeted Advertising Is Ruining the Internet and Breaking the World. *Motherboard*. https://motherboard.vice.com/en_us/article/xwjden/targeted-advertising-is-ruining-the-internet-and-breaking-the-world

Lapowsky, I., & Thompson, N. (2019, March 6). Facebook's Pivot to Privacy Is Missing Something Crucial. *Wired*. https://www.wired.com/story/facebook-zuckerberg-privacy-pivot/

Lehtiniemi, T. (2019). *Reorienting Datafication? New Roles For Users In Online Platform Markets | The Internet, Policy & Politics Conferences*. 2016. http://blogs.oii.ox.ac.uk/ipp-conference/2016/programme-2016/track-c-markets-and-labour/digital-markets-and-currencies/tuukka-lehtiniemi-reorienting.html

Lev-Aretz, Y. (2019, April). Facebook and the perils of a personalized choice architecture. *TechCrunch*.

http://social.techcrunch.com/2018/04/24/facebook-and-the-perils-of-a-personalized-choice-architecture/

Lindgren, S. (2019). Hacking Social Science for the Age of datafication. *Journal of Digital Social Research, 1(1), 1-9.*

Lyon, D. (2018). *The Culture of Surveillance: Watching as a Way of Life* (1 edition). Polity.

Maass, M., Wichmann, P., Pridöhl, H., & Herrmann, D. (2017). PrivacyScore: Improving Privacy and Security via Crowd-Sourced Benchmarks of Websites. *ArXiv:1705.05139 [Cs]*, *10518*, 178–191. https://doi.org/10.1007/978-3-319-67280-9_10

Mahieu, R. L. P., Asghari, H., & Eeten, M. van. (2018). Collectively exercising the right of access: Individual effort, societal effect. *Internet Policy Review*, 7(3). https://policyreview.info/articles/analysis/collectively-exercising-right-access-individual-effort-societal-effect

Mai, J.-E. (2016). Big data privacy: The datafication of personal information. *The Information Society*, *32*(3), 192–199. https://doi.org/10.1080/01972243.2016.1153010

Maiberg, E., & Grauer, Y. (2018, March 27). What Are "Data Brokers," and Why Are They Scooping Up Information About You? *Vice*. https://www.vice.com/en_us/article/bjpx3w/what-are-data-brokers-and-how-to-stop-my-private-data-collection

Maiberg, E., Koebler, J., & Cox, J. (2018, December 5). Internal Documents Show Facebook Has Never Deserved Our Trust or Our Data. *Motherboard*. https://motherboard.vice.com/en_us/article/7xyenz/internal-documents-show-facebook-has-never-deserved-our-trust-or-our-data

Mann, M., & Matzner, T. (2019). Challenging algorithmic profiling: The limits of data protection and anti-discrimination in responding to emergent discrimination. *Big Data & Society*, *6*(2), 2053951719895805. https://doi.org/10.1177/2053951719895805

Marres, N., & Gerlitz, C. (2016). Interface Methods: Renegotiating relations between digital social research, STS and sociology. *Sociological Review*, *64*, 21–46.

Marres, N. (2017). Digital Sociology. Cambridge: Polity Press.

Masson, E., van Es, K., Wieringa, M., (2020). Data Walking for Critical Data Studies: An Explorative Survey of Walking Methodologies. *Digital Culture & Education*, 11(1), 36-52

Matzner, T. (2016). Beyond data as representation: The performativity of Big Data in surveillance. *Surveillance & Society*, *14*(2), 197–210. https://doi.org/10.24908/ss.v14i2.5831

Milan, S. (2018). Digital Traces in Context| Political Agency, Digital Traces, and Bottom-Up Data Practices. *International Journal of Communication*, *12*(0), 21.

Milan, Stefania, & van der Velden, L. (2016). The Alternative Epistemologies of Data Activism. *Digital Culture & Society*, *2*(2), 57–74. http://dx.doi.org/10.25969/mediarep/991

Milan, Stefanija, & Agosti, C. (2019, February 7). *Personalisation algorithms and elections: Breaking free of the filter bubble*. Internet Policy Review. https://policyreview.info/articles/news/personalisation-algorithms-and-elections-breaking-free-filter-bubble/1385

Mittelstadt, B. (2016). Auditing for Transparency in Content Personalization Systems. *International Journal of Communication*, *10*(0), 12.

Munk, A. K. (2019, February 20). Introduction to semantic analysis with Cortext. Retrieved June 21, 2019, from Anders Kristian Munk website: https://medium.com/@AnthropologicalMachines/introduction-to-semantic-analysis-with-cortext-19f355b7289a

Neyland, D. (2016). Bearing Account-able Witness to the Ethical Algorithmic System. *Science, Technology, & Human Values*, *41*(1), 50–76. https://doi.org/10.1177/0162243915598056

Noble, S. U. (2018). *Algorithms of oppression: How search engines reinforce racism*. New York University Press.

Olejnik, L., Castelluccia, C., & Janc, A. (2014). On the uniqueness of Web browsing history patterns. *Annals of Telecommunications - Annales Des Télécommunications*, *69*(1), 63–74. https://doi.org/10.1007/s12243-013-0392-5

Papadopoulos, P., Kourtellis, N., & Markatos, E. P. (2018). *Cookie Synchronization: Everything You Always Wanted to Know But Were Afraid to Ask*. https://doi.org/10.1145/3308558.3313542

Pariser, E. (2012). *The Filter Bubble: What The Internet Is Hiding From You*. Penguin.

Pasquale, F. (2015). *The Black Box Society: The Secret Algorithms That Control Money and Information*. Harvard University Press.

Paßmann, J., & Boersma, A. (2017). Unknowing Algorithms. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 139–146). Amsterdam University Press. https://doi.org/10.5117/9789462981362

Pasternack, S. M. and A. (2019, March 2). *Here are the data brokers quietly buying and selling your personal information*. Fast Company. https://www.fastcompany.com/90310803/here-are-the-data-brokers-quietly-buying-and-selling-your-personal-information

Pierson, J. (2014). Interdisciplinary Perspective on Social Media, Privacy and Empowerment: The Role of Media and Communication Studies in Technological Privacy Research. *Digital Enlightenment Yearbook 2014*, 265–274.

Privacy International. (2019, April 16). *How Apps on Android Share Data with Facebook—Report*. Privacy International. http://privacyinternational.org/report/2647/how-apps-android-share-data-facebook-report

Raley, R. (2013). Dataveillance and Surveillance. In *"Raw Data" is an Oxymoron* (p. 192). The MIT Press. https://mitpress.mit.edu/books/raw-data-oxymoron

Reigeluth, T. B. (2014). Why data is not enough: Digital traces as control of self and self-control. *Surveillance & Society*, *12*(2), 243–254. https://doi.org/10.24908/ss.v12i2.4741

Rieder, B., & Röhle, T. (2017). Digital Methods. From Challenges to Bildung. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 109–124). Amsterdam University Press. https://doi.org/10.5117/9789462981362

Rocher, L., Hendrickx, J. M., & Montjoye, Y.-A. de. (2019). Estimating the success of re-identifications in incomplete datasets using generative models. *Nature Communications*, *10*(1), 3069. https://doi.org/10.1038/s41467-019-10933-3

Roesner, F., Kohno, T., & Wetherall, D. (2012). *Detecting and Defending Against Third-Party Tracking on the Web*. 155–168. https://www.usenix.org/conference/nsdi12/technical-sessions/presentation/roesner

Roesner, F., Rovillos, C., Saxena, A., & Kohno, T. (2019, June 6). *TrackingObserver: A Browser-Based Web Tracking Detection Platform*. https://trackingobserver.cs.washington.edu/

Rogers, R. (2013). *Digital Methods*. MIT Press.

Rogers, R. (2015). Digital Methods for Web Research. In *Emerging Trends in the Social and Behavioral Sciences* (pp. 1–22). American Cancer Society. https://doi.org/10.1002/9781118900772.etrds0076

Rogers, R. (2017). Foundations of digital methods: Query design. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 125–137). Amsterdam University Press. https://doi.org/10.5117/9789462981362

Rogers, R. (n.d.). *The Googlization Question, and the Inculpable Engine*.

Rogers, R. (2018). Social Media Research After the Fake News Debacle. *PARTECIPAZIONE E CONFLITTO*, *11*(2), 557-570–570. https://doi.org/10.1285/i20356609v11i2p557

Ryan, J. (2018, September 12). *Regulatory complaint concerning massive, web-wide data breach by Google and other "ad tech" companies under Europe's*

*GDPR*. Brave Browser. https://www.brave.com/blog/adtech-data-breach-complaint/

Ryan, J. (2018). *Behavioural advertising and personal data*. Brave. https://brave.com/Behavioural-advertising-and-personal-data.pdf

Sadowski, J. (2019). When data is capital: Datafication, accumulation, and extraction. *Big Data & Society*, *6*(1), 2053951718820549. https://doi.org/10.1177/2053951718820549

Sandvig, C., Hamilton, K., Karahalios, K., & Langbort, C. (2014, May 22). *Auditing Algorithms: Research Methods for Detecting Discrimination on Internet Platforms*. 64th Annual Meeting of the International Communication Association, Seattle, WA, USA. https://pdfs.semanticscholar.org/b722/7cbd34766655dea10d0437ab10df3a127396.pdf

Schmidt, C. D. (2018). *Google data collection research* (pp. 1–55). Digital Content Next, Vanderbilt University. https://digitalcontentnext.org/wp-content/uploads/2018/08/DCN-Google-Data-Collection-Paper.pdf

Schwartz, O. (2018, July 13). *Digital ads are starting to feel psychic*. The Outline. https://theoutline.com/post/5380/targeted-ad-creepy-surveillance-facebook-instagram-google-listening-not-alone

Seaver, N. (2013). Knowing algorithms. *Media in Transition 8*, 1–12. https://static1.squarespace.com/static/55eb004ee4b0518639d59d9b/t/55ece1bfe4b030b2e8302e1e/1441587647177/seaverMiT8.pdf

Seaver, N. (2017). Algorithms as culture: Some tactics for the ethnography of algorithmic systems. *Big Data & Society*, *4*(2), 2053951717738104. https://doi.org/10.1177/2053951717738104

*SOCIAL SCIENCE ONE*. (n.d.). Retrieved August 6, 2020, from https://socialscience.one/home

Store, P. D. (2016, December 27). *Facebook Ad Categories* [Text/html]. ProPublica Data Store. https://www.propublica.org/datastore/dataset/facebook-ad-categories

Sweeney, L. (2013). Discrimination in Online Ad Delivery. *COMMUNICATIONS OF THE ACM*, *56*(5), 44–54.

Szymielewicz, K. (2019, January 30). *Legal battle over online behavioural advertising widening*. Internet Policy Review. https://policyreview.info/articles/news/legal-battle-over-online-behavioural-advertising-widening/1384

Taylor, A., & Sadowski, J. (2015, May 27). *How Companies Turn Your Facebook Activity Into a Credit Score*. https://www.thenation.com/article/how-companies-turn-your-facebook-activity-credit-score/

Thompson, S. A. (2019, April 30). Opinion | These Ads Think They Know You. *The New York Times*. https://www.nytimes.com/interactive/2019/04/30/opinion/privacy-targeted-advertising.html, https://www.nytimes.com/interactive/2019/04/30/opinion/privacy-targeted-advertising.html

Tufekci, Z. (2019, March 8). Opinion | Zuckerberg's So-Called Shift Toward Privacy. *The New York Times*. https://www.nytimes.com/2019/03/07/opinion/zuckerberg-privacy-facebook.html

Uricchio, W. (2017). Data, Culture and the Ambivalence of Algorithms. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 125–137). Amsterdam University Press. https://doi.org/10.5117/9789462981362

van Dijck, José. (2013). 'You have one identity': Performing the self on Facebook and LinkedIn. *Media, Culture & Society*, 35(2), 199–215. https://doi.org/10.1177/0163443712468605

van Dijck, Jose. (2014). Datafication, dataism and dataveillance: Big Data between scientific paradigm and ideology. *Surveillance & Society*, 12(2), 197–208. https://doi.org/10.24908/ss.v12i2.4776

van Dijck, J., & Poell, T. (2013). Understanding Social Media Logic. *Media and Communication*, 1(1), 2–14. https://doi.org/10.17645/mac.v1i1.70

van Es, K., & Schäfer, M. T. (2017). Introduction. New Brave World. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 15–22). Amsterdam University Press. https://doi.org/10.5117/9789462981362

van Es, K., López Coombs, N., & Boeschoten, T. (2017). Towards a Reflexive Digital Data Analysis. In K. van Es & M. T. Schäfer (Eds.), *The Datafied Society. Studying Culture through Data* (pp. 171–180). Amsterdam University Press. https://doi.org/10.5117/9789462981362

van Es, K., Wieringa, M., and Schafer, M.T. (2018). Tool Criticism: From Digital Methods to Digital Methodology. In International Conference on Web Studies (WS.2 2018), October 3-5, 2018, Paris, France. ACM, New York, NY, USA.

Venkatadri, G., Andreou, A., Liu, Y., Mislove, A., Gummadi, K. P., Loiseau, P., & Goga, O. (2018). Privacy Risks with Facebook's PII-Based Targeting: Auditing a Data Broker's Advertising Interface. *2018 IEEE Symposium on Security and Privacy (SP)*, 89–107. https://doi.org/10.1109/SP.2018.00014

Venkatadri, G., Sapiezynski, P., Redmiles, E., Mislove, A., Goga, O., Mazurek, M., & P. Gummadi, K. (2019). Auditing Offline Data Brokers via Facebook's Advertising Platform. *Proceedings of the 2019*

*World Wide Web Conference (WWW'19)*, 1920–1930.
https://doi.org/10.1145/3308558.3313666

Venkatadri, Giridhari, Lucherini, E., Sapiezynski, P., & Mislove, A. (2019).
Investigating sources of PII used in Facebook's targeted advertising.
*Proceedings on Privacy Enhancing Technologies*, *2019*(1), 227–244.
https://doi.org/10.2478/popets-2019-0013

Venturini, T., Meunier, A (2019). Drafting and atlas on Artificial
Intelligence's matters of reflection. Available at:
http://www.tommasoventurini.it/ai/ . Last accessed January 29, 2020.

Venturini, T., Rogers, R. (2019). "'API-Based Research' or How Can Digital
Sociology and Digital Journalism Studies Learn from the Cambridge
Analytica Affair." Digital Journalism, Forthcoming.

Venturini, T., Bounegru, L., Gray, J., & Rogers, R. (2018). A reality
check(list) for digital methods. *New Media & Society*, *20*(11), 4195–
4217. https://doi.org/10.1177/1461444818769236

Wachter, S. (2019). *Affinity Profiling and Discrimination by Association in
Online Behavioural Advertising* [SSRN Scholarly Paper]. Social Science
Research Network. https://papers.ssrn.com/abstract=3388639

Wachter, S., & Mittelstadt, B. (2018). *A Right to Reasonable Inferences: Re-
Thinking Data Protection Law in the Age of Big Data and AI* (SSRN
Scholarly Paper No. ID 3248829). Retrieved from Social Science
Research Network website: https://papers.ssrn.com/
abstract=3248829

Weltevrede, E. J. T. (2016). *Repurposing digital methods: The research
affordances of platforms and engines* [University of Amsterdam].
https://dare.uva.nl/search?identifier=aaaa9bb3-8647-41df-954c-
2bb1e9f15d77

WhoTracksMe. (2019, July 18). *Trackers Rank*.
https://whotracks.me/trackers.html

Xu, Y. (Joe). (2018). Programmatic Dreams: Technographic Inquiry into
Censorship of Chinese Chatbots. *Social Media + Society*, *4*(4),
2056305118808780. https://doi.org/10.1177/2056305118808780

Zimmer, M. (2008). *The Gaze of the Perfect Search Engine: Google as an
Infrastructure of Dataveillance* (A. Spink & M. Zimmer, Eds.; pp. 77–
99). https://doi.org/10.1007/978-3-540-75829-7_6

Zuboff, S. (2015). Big other: Surveillance Capitalism and the Prospects of
an Information Civilization. *Journal of Information Technology*, *30*(1),
75–89. https://doi.org/10.1057/jit.2015.5

Zuboff, S. (2016, March 5). Google as a Fortune Teller: The Secrets of
Surveillance Capitalism. *Frankfurter Allgemeine Zeitung*.
https://www.faz.net/1.4103616

Zuboff, S. (2019). *The Age of Surveillance Capitalism: The Fight for a Human
Future at the New Frontier of Power*. Profile Books.