

Gillian Bolsover¹ and Philip Howard²

Computational propaganda has recently exploded into public consciousness. The US Presidential Campaign of 2016 was marred by evidence, which continues to emerge, of targeted political propaganda and the use of bots to distribute political messages on social media. This computational propaganda is both a social and technical phenomenon. Technical knowledge is necessary to work with the massive databases used for audience targeting; it is necessary to create the bots and algorithms that distribute propaganda; it is necessary to monitor and evaluate the results of these efforts in agile campaigning. Thus, a technical knowledge comparable with those who create and distribute this propaganda is necessary to investigate the phenomenon.

However, viewing computational propaganda only from a technical perspective—as a set of variables, models, codes and algorithms—plays into the hands of those who create it, the platforms that serve it and the firms that profit from it. The very act of making something technical and impartial makes it seem inevitable and unbiased. This undermines the opportunities to argue for change in the social value and meaning of this content and the structures in which it exists. Big data research is necessary to understand the socio-technical issue of computational propaganda and the influence of technology in politics. However, big data researchers must maintain a critical stance toward the data being used and analysed so as to ensure that we are critiquing as we go about describing, predicting or recommending changes. If research studies of computational propaganda and political big data do not engage with the forms of power and knowledge that produce it, then the very possibility for improving the role of social media platforms in public life evaporates.

Definitionally, computational propaganda has two important parts: the technical and the social. Focusing on the technical, Woolley and Howard define computational propaganda as the assemblage of social media platforms, autonomous agents, and big data tasked with the manipulation of public opinion³. In contrast, the social definition of computational propaganda derives from the definition of propaganda - communications that deliberately misrepresent symbols, appealing to emotions and prejudices and bypassing rational thought, to achieve a specific goal of its creators – with computational propaganda understood as propaganda created or disseminated using computational (technical) means.

Propaganda has a long history. Scholars who study propaganda as an offline or historical phenomenon have long been split over whether the existence of propaganda is necessarily detrimental to the functioning of democracies. However, the rise of the Internet and, in particular, social media has profoundly changed the landscape of propaganda. It has opened the creation and dissemination of propaganda messages, which were once the province of states and large institutions, to a wide variety of individuals and groups. It has allowed cross-border computational propaganda and interference in domestic political processes by foreign states. The anonymity of the Internet has allowed state-produced propaganda to be presented as if it were not produced by state actors. The Internet has also provided new affordances for the efficient dissemination of propaganda, through the manipulation of the algorithms and processes that govern online information and through audience targeting based on big data analytics. The social effects of the changing nature of propaganda are only just beginning to be understood and the advancement of this understanding is complicated by the unprecedented marrying of the social and the technical that the Internet age has enabled.

The articles in this special issue showcase the state of the art in the use of big data in the study of computational propaganda and the influence of social media on politics. This rapidly emerging field represents

¹ Postdoctoral Researcher in the Social Sciences, Oxford Internet Institute, University of Oxford

² Professor of Internet Studies and Director of Research, Oxford Internet Institute, University of Oxford

³ Woolley, Samuel C., and Philip N. Howard. "Automation, Algorithms, and Politics | Political Communication, Computational Propaganda, and Autonomous Agents — Introduction." *International Journal of Communication* 10, no. 0 (October 12, 2016): 9.

a new clash of the highly social and highly technical in both practice and research. We were brought on as guest editors of this special edition of Big Data to produce a more social science-focused edition of the journal. The process of reviewing the fifteen submissions, engaging with peer reviewers with both technical and social expertise, and closely editing the six papers published here has allowed us to reflect on the current status of this research area and offer suggestions for the future direction of the field.

Prediction, models and technical solutions should not be the primary goal of political big data research.

Almost all submissions to this special issue used big data to predict in some way—using social media data to predict levels of automation, outcomes of elections, or public opinion during referenda. Prediction was often seen to be a justification in and of itself. However, this should not be the case. We must evaluate the net academic contribution and social impact of predictive models, and be cognizant of the potential opportunities and costs of publicizing any predictive powers we develop as researchers. In short, big data can be immensely useful for making political inferences. However, developing the craft of prediction means improving the ability of many kinds of political actors to make political inferences. Solving social problems, redressing inequality and improving civic engagement are the kinds of outcomes we should strive for when we do our work.

Available doesn't mean ethical. Although the data of many social media sites are public, research has shown that users do not necessarily understand their information as such or that it could be used by researchers, companies or states. The same is true of consumer and other databases. More ethical questions arise when datasets are combined. Prior to the recent rise of socio-technical research areas, researchers in disciplines such as computer science, physics and engineering have rarely had to engage with the ethics of use of social data. It is important for social scientists to increase their technical expertise as politics moves online but also for technical fields to adhere to the professional ethics of social science research. Productive knowledge sharing and collaboration between previously-social and previously-technical disciplines are necessary so that big data studies can take the socially grounded and critical perspectives necessary for the study of social phenomena.

Don't throw in all the available data. If too many data features are included in a particular model, some spurious but apparently statistically-significant associations will arise. 'P-hacking' can occur even without deliberate malpractice: the problem of multiple comparisons means that even with stringent p-value cutoffs, proposing a sufficiently large number of models will always lead to false-positive inferences. This has led to reproducibility problems across the quantitative social sciences, and necessitates a re-evaluation of the best practices in our field. It is important to have a causal justification for all of the variables put into a model using big data approaches - we cannot simply put everything in the model and see what sticks. It is not only more efficient to begin with a small subset of variables for which there is theoretical support and contextual relevance and build a model methodically from there, but it is also absolutely essential if we want to avoid the danger of over-fitting false models.

Variables and models are important for what they tell us about underlying social phenomenon. In the same way as it is important to have a theoretical justification for analytical inputs, outputs must be evaluated for the knowledge they offer about underlying social phenomenon. Too many big data studies report only the predictive power of their models. However, prediction is not the goal. Understanding is the goal. Thus, each variable put into a model (which was based on a hypothesis derived from existing literature and understanding as to why it might be important) should be evaluated as to whether it was important or not in the models and what new knowledge this importance or lack thereof generates about the underlying social phenomenon being studied.

It is critically important to think about how research we produce might be used. One submission to this issue used Twitter data to attempt to predict whether protests would emerge based on social media data, arguing that although most protests are legal they cause disruption and damage to property and therefore it is important to predict them. The capability of predicting phenomenon, such as protests, crimes, elections and resignations, before they happen easily evokes a dystopian future in which this system of prediction is open to abuse and the potential for false positives that would undermined human agency and fundamental human rights. Big data researchers must not be complicit in making such dystopias a reality. The kinds of knowledge generated by technical studies of political big data are not simply truths; they are technologies and tools. We

must consider, in the work we put forward, whether the knowledge and tools we produce might be empowering those whose means and ends we would not wish to support.

Big data relies on what is available and obscures that which is not. An obvious limitation of big data research is that it relies on what is available. The majority of technical studies focus on Twitter data because the platform provides more open access than the more widely used Facebook. When a researcher queries a tweet using the Twitter API, the poster's join date, number of friends, number of followers and number of posts is returned. It is very common to see the analyses of these variables reported in papers because these data are available. Geographic location, religious affiliation, political preference, gender, level of education, and other variables commonly found to be strongly associated with social behaviors are almost impossible to access using Twitter data and, thus, much less is known about how such factors affect the concentration and circulation of computational propaganda. Although big data studies must rely on the data that is available, those who make use of this data must think critically about its political economy: Why has it been made available? Why was some data collected or made available and not others? Whose purposes does the existence of these data serve? What populations and issues are excluded from the dataset? In particular, we must not be satisfied with constructing studies around the data that are available but rather first decide what data is necessary to answer the research question and then seek to obtain it.

Focusing only on the technical prevents researchers from engaging in the social and, thus, opportunities for change. Big data has a great deal of value in social science research and there are newly arising opportunities for much greater understanding and ground-breaking research at the intersection of the technical and social. Computational propaganda is one of the current issues that this combination of the social and technical is necessary in order to understand. However, there is a danger in the focus on the technical that an understanding of the social is obscured. The conditions of production and reading are forgotten. Data is interpreted out of context. Research remains embedded in the structures of the current system and thus unable to engage with opportunities to change the current system. A more critical focus is necessary in big data studies of social phenomenon such that this research critiques the system rather than accepts the boundaries and structures of the system.

There is great opportunity in this emerging field. However, these opportunities must be based on collaboration and connection between knowledge and techniques derived from social and technical fields. The papers in this special issue sit more on the technical side of this research field but in compiling this special issue, we were careful to select only papers that incorporated social perspectives and contributed to social understanding.

In the first paper, Grimme, Preuss, Adam and Trautmann examine the issue of hybrid social bots, also sometimes referred to as cyborgs, that combine automation with human curation. After putting forward a definition and taxonomy of social bots, the authors present the results of an experiment in which they show that hybrid social bots are efficient for distributing messages, cost effective, and difficult to detect using automated means. Hybrid social bots are examples of phenomena that must be studied with an analytical frame that includes both the social and the technical.

The second and third papers focus on bot detection in understudied social contexts. In the second paper, Schäfer, Evert and Heinrich argue that bots and right-wing Internet activism comprised a semi-public sphere on social media in Japan's 2014 General Election. Grounding their data in a rich social context, they argue that incumbent President Shinzō Abe managed to appeal both to centrists, with his public pronouncements focusing on economics, and right-wingers, based on a hidden online nationalist agenda, and that this dual constituency was responsible for his electoral success.

In the third paper, Stukal, Sanovich, Bonneau and Tucker focus their attention on detecting bots in Russia's twittersphere, proposing a method that focuses on account properties and allows for retrospective analysis. They find that on the majority of days more than half of tweets using Russian political hashtags were produced by bots, shining a light on the functioning of computational propaganda in understudied non-democratic contexts.

The fourth and fifth papers move away from computational propaganda to consider political big data, more generally, and on using data to predict political outcomes. In the fourth paper, Sathiaraj, Cassidy and Rohli

combine voting data with consumer databases to attempt to predict election outcomes, using the US state of Louisiana as a case study. The way in which big data has been leveraged by modern political campaigns is a major issue that is rapidly evolving and this paper provides insight into how big data analysis can contribute to predicting – and implicitly following from this – shaping election outcomes.

In the fifth paper, Nigam, Dambanemuya, Joshi and Chawla turn the objects of big data prediction to a less studied phenomenon, peace processes, focusing on Twitter data collected around the 2016 Colombian referendum on a negotiated peace agreement. The authors argue that monitoring online data would have prevented the surprise rejection of the peace agreement in the referendum, but, more importantly, like many other papers in this issue they draw attention to the fact that the opportunities, pitfalls and effects of political big data can be very different in non-Western and non-democratic contexts.

The sixth and final paper in this special issue, by Huckle and White, turns attention back to the modern phenomenon of computational propaganda. Grounding their proposal in discussions of both the social and the technical, the authors propose implementing a mechanism based on blockchain technology to verify the provenance of news images. A great deal of attention has been paid to textual misinformation, but recent studies have shown how both photos and videos can be altered to spread propaganda messages. The proposals of this paper, although in their infancy, demonstrate how important it is for researchers of computational propaganda to advance their technical expertise to keep up with those who create and spread this propaganda.

Looking across the collection, conclude with several observations about the evolution of research into computational propaganda and political big data. First, there is increasing diversity in the range of political and institutional processes being treated for big data analysis. Elections, referenda, peace processes, and a full range of subnational, pre-election and post-election opinion formation processes are all active domains of inquiry. Second, standards for evidentiary quality are rising. Instead of casual dabbles in interesting social phenomena, the gold standard of research now involves the analysis of diverse political cultures in their own language and with enough author expertise to make the end product a credible piece of research to the area studies experts who know the culture, not just the methodologists who find big data work an intellectual puzzle. Third, the ethical challenges are getting more complex, not less. As guest editors of this special issue we struggled with a variety of ethical questions that led us to formulate the precepts for critical research agenda presented in this introduction.

We are proud to have been able to further this collaboration between the social and the technical through the production of this special edition of Big Data. We hope that this will be the beginning rather than the end of a merging of technically sophisticated methods and critical social perspectives about the influence of technology and big data in politics. Technology has shaken the foundations of established democracies and empowered extremists and authoritarians. It is imperative that the research community comes together, sharing knowledge and expertise across disciplines, to address the challenges of the merging of the social and the technical that have arisen in the Internet age.