

#### A SMART WEB FOR A MORE EQUAL FUTURE

A series focused on identifying the challenges and opportunities ahead and ways to address them

## ALGORITHMIC ACCOUNTABILITY

Applying the concept to different country contexts

July 2017

www.webfoundation.org

## CONTENTS

		· ////
For	eword	3
Introduction		4
01	The Opportunities	6
02	The Challenges	8
03	The Solutions	10
04	Case Example: Social Platforms	14
05	A Way Forward	16
06	References	18





The Web Foundation was established in 2009 by Sir Tim Berners-Lee, inventor of the World Wide Web. Our mission is to establish the open web as a public good and a basic right.

This paper has been adapted by the Web Foundation from a draft report commissioned to David Sangokoya of Data-Pop Alliance.

This research was funded by a grant from the Ford Foundation.

Copyright, World Wide Web Foundation, CC BY 4.0

### FOREWORD

"To achieve this vision, we must keep an eye on the trends, technologies and forces shaping the web of tomorrow, and the policy interventions that will be required to ensure digital equality becomes a reality." Welcome to this new series of policy white papers, produced by the World Wide Web Foundation.

The Web Foundation was established in 2009 by Sir Tim Berners-Lee, inventor of the World Wide Web. Our mission is to establish the open web as a public good and a basic right. Our five-year strategy – developed in 2016 – is to deliver digital equality – a world where everyone has the same rights and opportunities online. To achieve this vision, we must keep an eye on the trends, technologies and forces shaping the web of tomorrow, and the policy interventions that will be required to ensure digital equality becomes a reality.

On the web's 28th birthday in March 2017, Sir Tim Berners-Lee penned a letter on what he believed to be the biggest challenges facing the web today. The challenges he outlined are threefold: we've lost control over our personal data; misinformation spreads too easily online; and we need more transparency and understanding of digital political advertising.

Since then we have been discussing ways in which we could and should tackle these issues. We understood that these could be early warning signals of deeper problems, and set out to distil these in search of their most basic components. We landed upon data, algorithms and artificial intelligence, and the way these interact with existing socio-legal frameworks. These three issues are interdependent – data feed algorithms that are increasingly being used to make critical decisions, algorithms are the bedrock of artificial intelligence, and data gathered by AI and algorithms feed back into the system.

This is one of the three white papers we commissioned to begin to understand more about these issues. All too often, research, debate and discussion on these areas is focused on the US, UK and Europe, while actors from outside these countries are seldom being included as critical actors in thinking through policies at the global level. Our objective was to gain initial insights how each component is currently playing out in low and middle-income countries, and what some of the future risks and opportunities are.

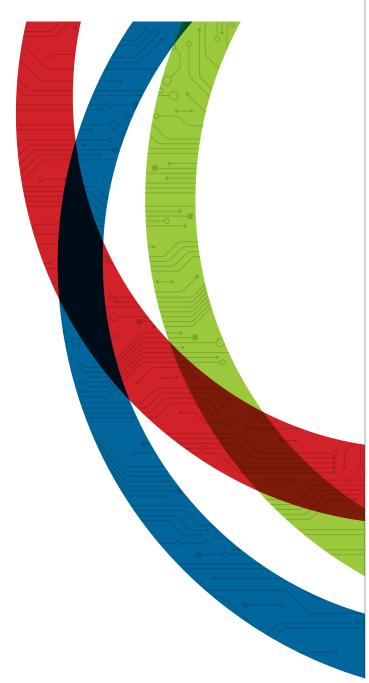
An important step towards enabling collaboration and solving the challenges the web faces is increasing public and key stakeholder understanding of how the individual components of the system work. We hope that these papers make a small contribution towards this goal, including in countries too often ignored in these debates. We will now be using these papers to refine our thinking and set our work agenda in the years ahead. We are sharing them openly in the hopes that they benefit others working towards our goals.

We hope you enjoy the read, and we welcome your feedback. Let's work together to build a more open web for a more equal world.

#### Craig Fagan

Director of Policy, Web Foundation June 2017

## INTRODUCTION



A t the centre of our information societies is the production of massive amounts of data through platforms, social networks, and machines. Data has not just become big; it is also increasingly becoming open and linked. In order to make sense of these huge amounts of data, and ensure their full richness is leveraged, companies and public sector actors are relying on algorithms — typically defined as structured set of rules for problem solving, executed by computers<sup>1</sup>. As such, algorithms are critical enablers of the data revolution that is taking place.

In the private sector, algorithms have become the backbone of many business models deployed worldwide. In the public sector — particularly in Europe and the US — algorithmic decision-making has emerged alongside broader policy trends of the last decade such as open government and evidence-based decision-making, and is now starting to be used in high-stakes areas such as criminal justice.

In low and middle-income countries, some governments and companies have more recently begun using algorithms in a development and public policy context. Recognising some longstanding failures in these sectors, governments and companies in these countries are now turning towards algorithmic systems as a way of balancing efficiency, fairness, and accountability in their decision-making processes.

As more tasks and decisions are delegated to algorithms, and they are provided more liberties in the way they execute such tasks, there is a growing concern: algorithms are controlling the inclusion — and exclusion — of people and information in an increasing number of settings. This grants algorithms the power to perpetuate, reinforce or even create new forms of injustice. Yet the outcomes of algorithmic processes are often not designed to be accessible, verified or evaluated by humans, limiting our ability to identify if, when, where, and why the algorithm produced harm — and worse still — redress this harm.

Who should be held accountable for the impact of algorithms, and what meaningful mechanisms — technical, legal, and policyoriented — should governments, companies, citizens and other stakeholders turn to for solutions?

David Sangokoya of Data-Pop Alliance provided the conceptual framing and draft for this paper, as part of their collaboration with the Web Foundation.

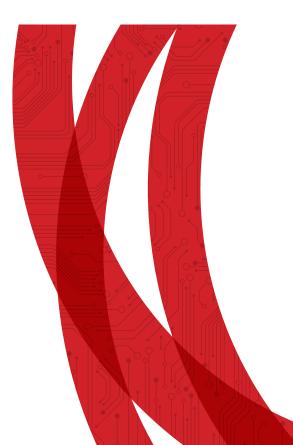
1 To avoid confusion, this paper employs a systems definition of "algorithm," describing code and data as well as the greater "socio-technical assemblage that includes algorithm, model, target goal, data, training data, application, hardware — and connect it all to a broader social endeavor." This paper argues that separating algorithmic accountability — the responsibility of algorithm designers to provide evidence of potential or realised harms — from algorithmic justice — the ability to provide redress for harms — is critical. The underlying reasons of algorithmic harms often lie in much larger and fundamental systemic issues.

Until now, much of the debate on how to accurately identify potential algorithmic bias and harms has occurred either within internal corporate research labs or within the academic research world, and there has been a lack of consensus amongst the broader community regarding what a "solutions toolkit" would look like.

Drawing from interviews with global experts, topic workshops and content research, this scoping paper aims to provide the reader with an understanding of algorithmic decision-making processes and the challenges they pose to our existing understanding of accountability across different contexts. It offers a map of existing technical and governance mechanisms for both identifying and addressing algorithmic harms and bias, as well as a set of recommendations and entry points for the Web Foundation and other stakeholders to contribute to this emerging field most effectively.



## THE OPPORTUNITIES



A the centre of our information societies is the production of massive amounts of data through platforms, social networks, and machines. Increasingly, companies have turned to automated machines and agents to make sense of this abundance of data through algorithms.

Although typically defined as a set of "encoded procedures" or "a logical series of steps for organising and acting on a body of data to quickly achieve a desired outcome"<sup>2</sup>, the term algorithm is often intended to describe a larger intersection of code, data and automated decisions. Originating from computer science and used in various social science disciplines, the term has been used to convey various meanings on the intertwining of human and machine decision inputs, and the extent to which the term includes code, data and ecosystems often varies.

Algorithms are growing in diversity and application as governments shift towards evidence-based decision-making. With mountains of data waiting to be mined, and algorithms' powerful ability to make statistical predictions and recommendations, it is no surprise that public sector actors are turning to algorithms to solve complex problems at the limits of human decision-making.

Algorithmic approaches and systems allow for collecting, classifying, structuring, aggregating and analysing data in such a way that unexpected insights, trends and predictions often become apparent (see Table 1). The history of human decision-making is wrought with examples of inefficient and unjust outcomes. The turn towards algorithms in governments — particularly in sectors such as criminal justice, healthcare, safety, fair employment and others — can be seen as part of a greater effort towards evidence-based decision-making and the adoption of open and transparent government principles.

2 Gillespie, T. (2014). The Relevance of Algorithms. Media technologies: Essays on communication, materiality, and society, 167.

Governments and companies working in low, and middle-income countries have also more recently begun deploying algorithms, largely in development and public policy context (Table 2, below, provides some examples).

Machine learning, a field that has exploded over the recent years, has added a new layer of complexity since the algorithms it unleashes are capable of learning implicit rules from the data they are exposed to. The training process can involve in practice adjusting for millions of parameters generating "astronomically more possible outcomes than any [non-machine learning] algorithm could ever hope to try."<sup>5</sup> Many of the examples presented in this whitepaper are cases of the deployment of these complex machine-learning algorithms.

#### Table 1 — Examples of algorithms by function<sup>3</sup>

FUNCTION	ТҮРЕ	EXAMPLES
PRIORITISATION:	General search engines	Google, Bing, Baidu
associating rank with emphasis on	Special search engines	Genealogy, image search, Shutterstock
particular information or results at	Meta search engines	Info.com
the expense of others through a	Questions & answers	Quora, Ask.com
set of pre-defined criteria	Social media timelines	Facebook, Twitter
CLASSIFICATION:	Reputation systems	Ebay, Uber, Airbnb
grouping information based on	News scoring	Reddit, Digg
features identified within the	Credit scoring	Credit Karma
source data	Social scoring	Klout
ASSOCIATION:	Predictive policing	PredPol,
determining relationships between particular entities via semantic and connotative abilities	Predicting developments and trends	ScoreAhit, Music Xray, Google Flu Trends
FILTERING:	Spam filter	Norton
including and/or excluding	Child protection filter	Net Nanny
information as a result of a set of	Recommender systems	Spotify, Netflix
criteria	News aggregators	Facebook News Feed

#### Table 2 — Examples of algorithms in low and middleincome countries<sup>4</sup>

CATEGORY OF USE	SECTORS	EXAMPLES
EXPERT-LED RESEARCH	Public health	Event and Pattern Detection Lab (India, Sri Lanka)
INITIATIVES AND TOOLS	Crime	CrimeRadar (Brazil)
	Policing	CMore (South Africa)
	Education	Microsoft-Andhra Pradesh State (India)
GOVERNMENT DECISION SUPPORT	Credit scoring	Tala Mobile (formerly Inventure) (Kenya) Branch (Kenya)
	Transport (Ridesharing)	99Taxis (Brazil) Ola cabs (India)
CURATED KNOWLEDGE DISCOVERY	Personalised ed-tech	Geekie (Brazil)

Adapted from Lepri, Bruno, Staiano, J., Sangokoya, D., Letouzé, E., & Oliver, N. (2017). "The Tyranny of Data? The Bright and Dark Sides of Data-Driven Decision-Making for Social Good." Transparent Data Mining for Big and Small Data. Springer International Publishing,

For further detail see Appendix II, and one of our other documents in this series: "Maximising the Opportunities and Minimising the Risks of Artificial Intelligence in Low and Middle Income Countries", available through our website http://webfoundation.org/ (Last accessed 5/2/2017) White House Executive Office of the President (2016, December), "Artificial Intelligence, Automation, and the Economy." https://www.whitehouse.gov/sites/whitehouse.gov/ files/images/EMBARGOED%20AI%20Economy%20Report.pdf (Last accessed 5/2/2017)



## THE CHALLENGES



Many algorithms either regulate the content we are presented with online, or rely on data that was extracted from the online space to provide offline products and services. Given that the internet is a global infrastructure providing access to a borderless space, there will be cases in which these algorithms will affect populations that are geographically and culturally distant to the places where the algorithms were designed. This increases the potential for algorithms to cause harm.

#### Algorithmic Harm

What do we mean by harm? Underlying the definition of harm are the values of a society. By defining what an algorithm should not do (harm), hard boundaries emerge for what an algorithms' optimisation function should be (broader objectives). Ensuring algorithms are compatible with the diversity of values worldwide is certainly a challenge. Who should define and determine whether there have been any harms produced by the algorithms? In what cases should we promote that those who might be affected by the algorithm be integrated into the design process?

This paper does not seek to provide a full-fledged definition of harm, but acknowledges the usefulness of a working definition. Algorithms are a relatively new component of the broader debate on harms. It therefore seems reasonable to borrow a definition of harm from the legal sphere, which defines harms as setbacks to interests that are also considered to be a "wrong" (i.e. something that is inflicted unfairly, not voluntarily consented to, or illegitimate).<sup>6</sup>

The issue of algorithmic harms has arisen primarily in public and expert discussions around the use of algorithms deployed by governments of high-income countries<sup>7</sup> in high-stakes sectors such as criminal justice, credit and employment. The harms here result from discrimination.

### Algorithmic Discrimination

Discrimination can occur in two ways. Two people may be the same in relevant aspects but are treated differently (such as two defendants committing the same crime, but one getting a lighter sentence). Or relevant differences between them are not accounted for, and the two people are treated in the same way (for example someone's zip code is used as one of the factors to determine the likelihood of defaulting on a loan). The failure to acknowledge these relevant details about an individual is what makes the outcome unfair, and therefore a wrong. In this way a person may reasonably expect for a certain outcome (like getting access to a well-paying job that they are qualified for) which is wrongfully prevented by an algorithm, constituting a harm.

Several authors and experts in the emerging field — such as Cathy O'Neil (Author, "Weapons of Math Destruction") and Julia Angwin (Senior Reporter, ProPublica) — have underlined the risks of algorithmic discrimination. These include numerous examples, such as the 2012 plans of a German credit agency to mine Facebook

<sup>6</sup> Feinberg, J. (1984). Moral limits of the criminal law. Volume 1, Harm to others (Oxford scholarship online). New York ; Oxford: Oxford University Press.

<sup>7</sup> Throughout this paper and all the papers that are part of this series we rely on the country classification system proposed by the World Bank. This system categorizes countries into 4 groups based on estimates of the per capita Gross National Income of the previous year: low-income (\$1,045 or less), middle-income (\$1,045 - \$4,125), upper-middle-income (\$4,125 - \$12,736), and high-income (12,736 or more). More information available through the WB website: World Bank Data Team (2015, July 7). New Country Classifications. The Data Blog https://blogs.worldbank.org/opendata/newcountry-classifications (Last accessed 5/2/2017)

and other social media data to assess creditworthiness,<sup>8</sup> and models to predict the probability that a given convict will reoffend.<sup>9</sup>

High and low income countries face the same categories of harms and threats from algorithmic decision-making. However, the impact of these harms can be vastly different depending on existing legal protections and accountability mechanisms in place, especially for marginalised groups. In some countries, algorithmic discrimination and inaccurate predictions may result in unwanted advertising or other inconveniences in customer experiences. But for marginalised groups in fragile contexts, however, many argue that algorithmic discrimination may lead to unchecked aggression, and even lifethreatening exclusion from public services and resources.

## The Causes of Algorithmic Discrimination

Algorithmic decision-making procedures can reproduce and reinforce existing patterns of discrimination, for example by inheriting the prejudice of prior decisionmakers, or by reflecting widespread biases that persist in society.<sup>10</sup>

Discrimination by algorithms can materialise as the result of problems at different stages.

- Biased or otherwise poor quality input data: The data may be biased, incomplete, or of otherwise poor quality, potentially leading an algorithm to produce poor and perhaps discriminatory outcomes.
  - For example, predictive policing often relies on previous arrests to define where police should be deployed, and the characteristics officials should search for in defining their targets. If we assume a specific minority is discriminated against, and thus its members are frisked more often than members of other groups, all things equal the arrests for possession of illegal drugs and undeclared weapons amongst members of this group should be reported disproportionately. If the data on arrests is not fed to the algorithm as a proportion of frisks, the algorithm would likely recommend continuing such disproportionate activity in those neighbourhoods and targeting members of the affected group.
- Poorly defined rules: The data used as an input for algorithmic decisions may be poorly weighted.

For example

Social credit scoring companies such as Kreditech and Tala Mobile, may lower a customer's credit limit not based on the customer's payment history, but rather based on location and social analyses of other customers with a poor repayment history that had shopped at the same establishments where the prospective customer had also shopped.<sup>11</sup>

- Overemphasis of zip code within predictive policing algorithms in the US can lead to the association of lowincome African-American neighbourhoods with areas of crime, and, as a result, the specific targeting based on group membership.<sup>12</sup>
- Lack of contextual awareness: The definition of quality of the training data and the robustness of the rules and weights is often context specific. Algorithms that work well within the context for which they were designed might discriminate if rolled out in a different context.<sup>13</sup>
  - For example, face detection software fails to detect the faces of minority groups, yet detects a face when shown a white mask.<sup>14</sup>
- **Feedback loops:** Algorithms don't operate in a vacuum. Their activity affects the environment from which they extract the data they use as input. A biased algorithm might reinforce its biases, in what could be deemed a self fulfilling prophecy loop.
  - For example, an algorithm might suggest (based on biased data or a glitch in the rules) that a specific group should be denied access to credit due to a perceived lack of capacity to repay. If this same algorithm is used widely enough the systematic exclusion of that group's access to credit might follow. Over time the economic well being of this whole group will deteriorate. The algorithm will have undermined members of this group individually, but also the informal social networks each individual member of the group relies on in moments of urgency. The group as such will have become, in effect, less capable of repaying loans.

All of the above four areas are interrelated. Algorithms can be seen as part of a broader system where data and their collection process, the rules that govern the algorithm, the decision-making process that follows, and the broader socio-legal frameworks are all interconnected. The question is: how does one hold to account an algorithm for such critical decisions?

<sup>8</sup> Medick, V. (2012, June 7) German Agency to Mine Facebook to Assess Creditworthiness. Der Spiegel, Available at http://www.spiegel.de/international/germany/germancredit-agency-plans-to-analyze-individual-facebook-pages-a-837539.html (Last accessed 5/2/2017)

<sup>9</sup> Rudin, C. (2015, September 9) New models to predict recidivism could provide better way to deter repeat crime. The Conversation. Available at http://theconversation.com/ new-models-to-predict-recidivism-could-provide-better-way-to-deter-repeat-crime-44165

<sup>10</sup> Crawford, K., Schultz, J. (2014) Big data and due process: Toward a framework to redress predictive privacy harms. Boston College Law Review 55(1), 93–128 ; Barocas, S., Selbst, A. (2016): Big data's disparate impact. California Law Review 104, 671–732

<sup>11</sup> Ramirez, E., Brill, J., Ohlhausen, M., McSweeney, T.: (2016) Big data: A tool for inclusion or exclusion? Tech. rep., Federal Trade Commission

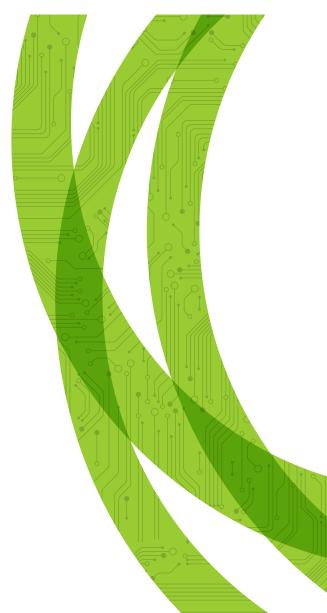
<sup>12</sup> Christin, A., Rosenblatt, A., boyd, d.: (2015) Courts and predictive algorithms. Data & Civil Rights Primer

<sup>13</sup> Calders, T., Zliobaite, I.: Why unbiased computational processes can lead to discriminative decision procedures. In: B. Custers, T. Calders, B. Schermer, T. Zarsky (eds.) Discrimination and Privacy in the Information Society, pp. 43–57 (2013)

Discrimination and Privacy in the information Society, pp. 43–57 (2013) 14 Algorithmic Justice League (n.d.) The Coded Gaze. Algorithmic Justice League Available at http://www.ajlunited.org/the-coded-gaze (Last accessed 5/2/2017)



## THE SOLUTIONS



#### Algorithmic Accountability

 ${f T}$ o begin to address algorithmic harms and discrimination, the concept of algorithmic accountability has begun to emerge.

Accountability is usually referred to as the duty governments and other authorities have to present themselves before those whose interest they represent or are otherwise bound to, and justify how power was exercised, and resources were used.<sup>15</sup>

When applied to algorithms, algorithmic accountability has often been conflated with other values, such as transparency.<sup>16</sup> Transparency has been held as an essential component of accountability, enabling citizens, consumers, data journalists, watchdog organisations and others to verify and understand the inputs, processes and outputs of a complex algorithmic system to identify evidence of harms as a first step for redress.<sup>17</sup>

However, several researchers in recent years have pointed to limitations in defining algorithmic accountability as transparency. Crawford and Ananny (2016) classified and filtered these into a list of 10 of transparency's limitations. The list includes the claim that the new complexities introduced by algorithms make "being able to see a system" as insufficient for "being able to know how it works and [to] govern it."<sup>18</sup>

15 Lister, S. (2010). Fostering Social Accountability: From Principle to Practice-A Guidance Note. New York: United Nations Development Programme; , McGee, R. and J. Gaventa, eds. (2010b). Review of Impact and Effectiveness of Transparency and Accountability Initiatives: Synthesis Report. Prepared for the Transparency and Accountability Initiative Workshop October 14–15. Open Society Institute.

- 16 Saurwein, F., Just, N., & Latzer, M. (2015). Governance of algorithms: options and limitations. info, 17(6), 35-49.
- 17 N. Diakopoulos (2014). Algorithmic Accountability Reporting: On the Investigation of Black Boxes. Tow Center.
- 18 Ananny, M., & Crawford, K. (2016, December). Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability. new media & society, 1461444816676645;; Janssan, M., and Kuk, G. (2016, July) "The Challenges and Limits of Big Data Algorithms in Technocratic Governance." Government Information Quarterly 33; Kroll, J.A., Huey, J., Barocas, S., Felten, E.W., Reidenberg, J.R., Robinson, D.G., Yu, H.: (2017, February). Accountable algorithms. University of Pennsylvania Law Review

Although we are at a stage in which the definition of algorithmic accountability is still being agreed upon, experts and practitioners have been putting forward general principles to be debated.

In January 2017 the Association for Computing Machinery (ACM) called for comments regarding its statement on the growing risk of algorithmic bias, where in defining the practical implications of accountability in this context, it claimed that "Institutions should be held responsible for decisions made by the algorithms that they use, even if it is not feasible to explain in detail how the algorithms produce their results."<sup>19</sup>

Another actor at forefront of interdisciplinary discussions on reevaluating how accountability in particular can be more clearly defined is the Fairness, Accountability and Transparency in Machine Learning (FATML) community. This interdisciplinary academic community of computer scientists, developers and researchers organised itself in 2014. In 2016 FATML released-and opened for comment-a set of five guiding principles for "accountable algorithms", which they hope will "help developers and product managers design and implement algorithmic systems in publicly accountable ways. Accountability in this context includes an obligation to report, explain, or justify algorithmic decision-making as well as mitigate any negative social impacts or potential harms."<sup>20</sup>

### **TABLE 3** — Principles for Accountable Algorithms – Fairness, Accountability and Transparency in Machine Learning (2016)

PRINCIPLE	DESCRIPTION
FAIRNESS	"Ensure that algorithmic decisions do not create discriminatory or unjust impacts when comparing across different demographics"
EXPLAINABILITY	"Ensure that algorithmic decisions as well as any data driving those decisions can be explained to end-users and other stakeholders in non-technical terms."
AUDITABILITY	"Enable interested third parties to probe, understand, and review the behaviour of the algorithm through disclosure of information that enables monitoring, checking, or criticism, including through provision of detailed documentation, technically suitable APIs, and permissive terms of use."
RESPONSIBILITY	"Make available externally visible avenues of redress for adverse individual or societal effects of an algorithmic decision system, and designate an internal role for the person who is responsible for the timely remedy of such issues."
ACCURACY	"Identify, log, and articulate sources of error and uncertainty throughout the algorithm and its data sources so that expected and worst case implications can be understood and inform mitigation procedures."

<sup>19</sup> Association for Computing Machinery (2017). Statement on Algorithmic Transparency and Accountability: Association for Computing Machinery US Public Policy Council (USACM)http://www.acm.org/binaries/content/assets/public-policy/2017\_usacm\_statement\_algorithms.pdf (Last accessed 5/2/2017) In December 2016, the Institute of Electrical and Electronics Engineers (IEEE) released a draft for public discussion titled "Ethically aligned design", where regarding the more autonomous type of algorithms used for ai it described accountability under the umbrella of responsibility, and underlined the need for accountability that can help "proving why a system operates in certain ways to address legal issues of culpability, and to avoid confusion or fear within the general public". See Ethically Aligned Design. IEE. Available at http://standards. ieee.org/develop/indconn/ec/ead\_v1.pdf (Last accessed 5/2/2017)

<sup>20</sup> Diakopoulos, N. et al. (n.d) Principles for Accountable Algorithms and a Social Impact Statement for Algorithms. Fairness, Accountability, and Transparency in Machine Learning (FATML) http://www.fatml.org/resources/principles-for-accountable-algorithms (Last accessed 5/2/2017)

Although such general principles shed light on important aspects of algorithmic accountability, it is essential for companies and governments to find new methods and options for accounting for harms. This can happen by focusing on what conditions or categories of harms companies and governments should be accountable for, and to whom they should be accountable to.<sup>21</sup> While citizens are often the focal point of accountability efforts, instances involving sensitive information may not allow for direct citizen engagement. Accounting in these circumstances may occur before government oversight committees, internal auditors or regulators who would be granted greater access to the details of the system and responsibilities to investigate specific elements (such as disparate impact).

While companies and governments have the responsibility to account for harms related to algorithmic decision-making, it is not definitively clear who should have the responsibility to repair such harms. Take for example Facebook and the issue of fake news propagated with social media platforms. While Facebook has the responsibility to take account for fake news happening within its system, changing the incentive structures that drive fake news as well as redressing the harms it might generate is a much bigger task that will involve Facebook as well as other actors. This leads us to the concept of algorithmic justice.

#### Towards Algorithmic Justice: Clarity on Responsibilities

Algorithmic decision-making introduces new complexities to existing forms of accountability and redress of harms. The nature and complexity of most algorithms in question make them "black boxes" which limit the ability to query their processes and decisions, putting citizens and consumers at a high level of risk. Viable approaches to redress any resulting harms from algorithms could open the way for algorithmic justice. Understanding responsibility for the functioning of the various elements of algorithms is an important step on the path towards establishing justice.

### Designer responsibility: interpretability, oversight and dynamic testing

Designers of algorithms are the people or institutions that have established AN algorithm's rules, weights and/or inner workings. The category of designers can be taken by the same government units and private companies that are leveraging algorithms to process data, but often the task of designing algorithms falls under the responsibility of specialized third-party vendors.

Through organizations such as FATML and academic workshops, corporate and academic researchers have been exploring new methodologies and techniques for governments and companies to account for potential harms associated with algorithms. Accounting options are inherently technical, and aimed at testing existing algorithms for potential harms. Much of the conversation on how to identify areas of potential algorithmic harms occurs either within internal corporate research labs or within the academic research world, and there is a lack of consensus on a "solutions toolkit" with the exact methods and tools for accounting at a general level.

However there are three key areas relevant for governments and companies trying to find evidence of algorithmic harms that merits exploration: interoperability, oversight and dynamic testing.

- Interpretability: Governments and companies generally rely on out-of-the-box algorithmic products purchased or licensed through third-party vendors to deploy algorithmic systems, but are often not equipped to evaluate candidate vendors' offers beyond a simple cost comparison. In order to account for potential harms and account for specific risks prior to implementation, governments and companies need frameworks to evaluate how proposed systems function and perform with regard to multiple categories of harms.
- Oversight: With the closest access to algorithmic systems, governments and companies can employ forms of auditing and reporting either for their internal measures (or for regulators), through trusted intermediaries (such as universities or companies) or to the public. These kinds of audits will evaluate input data, decision factors, and output decisions and involve documentation and potentially API access and permissive terms of use. Forms of auditing and reporting could include: discretionary internal audits, periodic internal audits, third-party audits, regulatory audits, and public audits.

21 Interview with Anupam Datta (Carnegie Mellon University)

Dynamic testing: There needs to be further exploration on how computational methods could provide accountability within black box systems.<sup>22</sup> Dynamic testing could offer such an option by employing cryptographic commitments (i.e. equivalents of sealed documents held by a third party or in a safe place); fair random choices (i.e. a technique allowing software to make fully reproducible random choices); or zero-knowledge proofs (i.e. cryptographic tools that allow a decision-maker to prove that the decision policy that was actually used has a certain property without revealing either how the property is known or what the decision policy is).

#### Social responsibility: journalism and literacy

Although governments and companies retain the responsibility to account for harms, computational journalism and efforts towards algorithmic literacy in citizens can be essential for citizens and citizen-centric groups to understand their own participation in algorithmic systems.

- **Computational journalism:** Computational journalism has evolved over time from "the application of computing technology to enable journalism across information tasks" to the investigation of algorithmic systems in order to characterise their functions, power, and biases as "algorithmic accountability reporting." As described in Diakopoulos (2015), examples of computational journalism include "comparisons and visualizations of statistical models of unemployment correction, to sophisticated reverse engineering investigations of online price discrimination."<sup>223</sup>
- Algorithmic literacy: Algorithmic literacy involves efforts to "enable more individuals to impact information flows and perceive when or if they or others are being marginalized."
  <sup>24</sup> However, the impact of these efforts may be limited as a result of the level of technical knowledge needed, as well as the discrepancy in identifying manipulation and being able to make significant changes to the status quo. Additionally, algorithmic literacy efforts would need to be context and industry-specific, and assume limited changes in the algorithm over time.

Incentives to invest in further resources to account for algorithmic harms is currently only as powerful as public pressure and journalistic exposure, which manifests through public information (e.g. algorithmic accountability reporting, investigative research, etc.). In a recent article, "Towards Accountability: Data, Fairness, Algorithms, Consequences," Danah Boyd highlights that while companies may not know how the systems that they design will evolve, they also unfortunately don't "have the tools to know when they're being gamed and when their systems are being manipulated or used to do harm."<sup>25</sup>

Clearly, the conversation about how to establish algorithmic justice is only just beginning.

<sup>22</sup> Kroll et al. (2017, February).

<sup>23</sup> Diakopoulos, Nicholas. "Algorithmic accountability: Journalistic investigation of computational power structures." Digital Journalism 3.3 (2015): 398-415; Rainie, L. & Anderson, J.Q. (2017, February 8). Code-Dependent: Pros and Cons of the Algorithm Age. Pew Internet & American Life Project. Available at http://www.pewinternet. org/2017/02/08/code-dependent-pros-and-cons-of-the-algorithm-age/, (Last accessed 5/2/2017)

<sup>24</sup> Caplan, Robyn, and Reed, L. (2016, May 16) Who Controls the Public Sphere in an Era of Algorithms? Data & Society Research Institute. Available at https://datasociety.net/ pubs/ap/CaseStudies\_PublicSphere\_2016.pdf (Last accessed 5/2/2017)

<sup>25</sup> boyd, d.(2017, April 12) Towards Accountability: Data, Fairness, Algorithms, Consequences. Data & Society Research Institute. Available at https://points.datasociety.net/ toward-accountability-6096e38878f0 (Last accessed 5/2/2017)



## CASE EXAMPLE: SOCIAL PLATFORMS



Consumers in both high- and low-income countries constantly engage with algorithmic decision-making when they access the internet and receive curated content, particularly via social platforms. As more users come online in the coming years in low and middle-income countries, the impacts of this reality will only become amplified. Looking at high-income markets provides some clues for what may be in store.

Social media has increasingly become a major source and distributor of news and information for citizens in high-income countries. A Pew Research study conducted in late 2016 on digital news in the US highlighted that users were "equally likely to get news by going directly to a news website (36 per cent) as getting it through social media (35 per cent)."<sup>26</sup> Additionally, the same study indicated that 10 per cent of consumers cite 'Facebook' as a specific news outlet.

The current business model for social media companies promotes a like-minded "information concentration" on their platforms. Companies gain advertising profit from increased user engagement and interactions on their platforms, and the algorithms behind search engines and social media platforms are therefore designed to manipulate our online experiences towards information users like (i.e. "filter bubbles"), amplifying biases and distorting perspectives.<sup>27</sup> Additionally, disinformation groups in the U.S. have high incentives to specifically target Facebook communities<sup>28</sup> to earn thousands of dollars generated advertising revenue from clicks and visits on their websites. This phenomenon has been written about widely, particularly in the context of how misinformation spread on social media may have influenced major election and referendum results in high-income countries.

Yet the phenomenon is occurring in high and low-income countries alike. Misinformation, in the form of "fake news," disinformation (e.g. hoaxes) and propaganda are not new, but the speed at which false news stories are distributed across algorithmically curated social media platforms has negatively impacted user perception and interpretation of factual information.

**India:** With over 160M monthly users (compared to Facebook's 155M users), Whatsapp has been the primary vehicle for fake news and rumour spreading in India, with deleterious effects. In 2015, rumours of gang violence quickly led to the formation of new gangs, and incorrect information on government shortage of salt led to abrupt panic in grocery shops and the trampling of a woman.<sup>29</sup> In 2016 an elaborate hoax article detailing the insertion of traceable nano GPS chips into new government bank notes led to riots and communal violence before eventually being debunked.<sup>30</sup>

**Indonesia:** In Indonesia's 2014 presidential election (in which political and religious group affiliation claims have previously inflamed deadly violence), smear campaigns and false images depicted President Joko Widodo, a Muslim moderate, as a Christian of Chinese descent. His party was forced to quell the rumours by posting a photo of his identity certificate on Facebook.<sup>31</sup>

27 Diakopoulos, N. and Frielder, S. (2016, November) ; Pariser 2011

<sup>26</sup> Mitchell, A. (2017, February 9) How Americans Encounter, Recall and Act Upon Digital News. Pew Research Center http://www.journalism.org/2017/02/09/ how-americans-encounter-recall-and-act-upon-digital-news/ (Last accessed 5/2/2017)

<sup>28</sup> Dredge, S (2015, July 15). Facebook and Twitter on the rise as sources of news in the US, The Guardian https://www.theguardian.com/technology/2015/ jul/15/facebook-twitter-sources-news-pew (Last accessed 5/2/2017)

Smith, O. (2013, February 13) WhatsApp fake news crisis is leading to riots & bloodshed. The Memo http://www.thememo.com/2017/02/13/whatsappindia-fake-news-crisis-is-leading-to-riots-bloodshed/ (Last accessed 5/2/2017)
 Ibid

<sup>31</sup> Kwok, Y. (2017, January 5) Where Memes Could Kill: Indonesia's Worsening Problem of Fake News. Time. Available at http://time.com/4620419/indonesiafake-news-ahok-chinese-christian-islam/ (Last accessed 5/2/2017)

**Myanmar:** Recent telecommunication reforms following Myanmar's decades under military regimes gave internet access for the first time to over 51 million Burmese people, "leapfrogging the era of dial-up and desktops [and] starting with mobile phones and social media."<sup>32</sup> Persisting social tension between majority Buddhist and minority Muslim population tinge these new online exchanges, with new "Facebook-first" media organisations pushing extremist propaganda.<sup>33</sup>

**South Africa:** During the 2016 municipal election in South Africa, a satirical website published a fake report claiming that one of the political parties had illegally marked almost 100,000 ballots ahead of the election.<sup>34</sup> Although the Independent Electoral Commission debunked the report, the report had already been circulated and read by tens of thousands of voters. Although the website was deleted in violation of the country's Municipal Electoral Act, the content remains available and searchable across cached versions and shares of the original content in news index.

Popular false news stories may spread differently in low and middle-income countries as compared to other countries due to:

- Limited information on government statistics and open data<sup>35</sup> for baseline and further investigation.
- Limited independent news media sources<sup>36</sup> to foster journalistic ethics, provide trusted and unbiased coverage of events separate from government.
- **Higher relative cost of internet access**,<sup>37</sup> which limit the amount of time people might be willing to spend online looking for, or encountering, counterarguments.
- **Higher prominence of mobile-only internet access**, which limits the type of activities users engage in, increasing the passive aspects of information consumption, over more active forms of engagement with information.<sup>38</sup>
- **Language barriers**<sup>39</sup> **and higher illiteracy rates**,<sup>40</sup> which limit the universe of resources effectively accessible to users over the internet (national or foreign), which might provide robust counter-arguments to the untruthful piece of information that is being spread.

Currently, there are very few common standards or ground rules that apply for social media news distribution. In response to recent critiques, Facebook has taken steps towards making instances of fake news more interpretable (e.g. displaying original new sources in headlines, giving notification to users of questionable news articles prior to user sharing, etc.) and providing limited forms of auditing (in collaboration with trusted intermediaries and partnerships) to identify news content flagged by several auditors as questionable.<sup>41</sup>

Yet experts note that the inherent problems are cross-sector and not purely technical,<sup>42</sup> which require a multi-stakeholder approach:

- **From journalism, news media and civic tech:** verification and standard tools for fact-checking (e.g. Check).
- From legal and regulatory authorities: monitoring legislation that respects freedom of expression and speech. The recent German and Nigerian examples of a regulatory response are ones that may be well-intentioned but some have suggested they are misguided and could ultimately lead to self-imposed censorship. The German government has discussed a plan to "force social media companies to monitor and censor some kinds of online expression."<sup>43</sup> The Nigerian government also aims to adapt regulations on hate speech and libel in order to address concerns around online misinformation.<sup>44</sup>
- From the public: algorithmic literacy trainings and guides; individual browser extensions to embed fact-checking scripts.
- From social media platform companies: design transparency (e.g. decision factors for flagging certain content); partnerships to identify problematic URLs across platforms.

38 Napoli, P., & Obar, J. (2014). The Emerging Mobile Internet Underclass: A Critique of Mobile Internet Access. The Information Society, 30(5), 323-334.

<sup>32</sup> Connoly, K. (2016, December 2) Fake news: an insidious trend that's fast becoming a global problem . the Guardian, Available at https://www.theguardian.com/media/2016/ dec/02/fake-news-facebook-us-election-around-the-world (Last accessed 5/2/2017)

<sup>33</sup> Connoly, K., ibid.

<sup>34</sup> De Wet, P ( 2016, August 5) Fake news websites fall foul of the IEC after marked ballot paper story. Mail and Guardian. Available at https://mg.co.za/article/2016-08-05-00fake-news-websites-fall-foul-of-the-iec-after-marked-ballot-story-earlier-this-week (Last accessed 5/2/2017)

<sup>35</sup> Open Data Barometer, 4th Edition (2017). Web Foundation. Available at http://opendatabarometer.org/ (Last accessed 6/2/2017)

<sup>36</sup> Freedom House (2017) Press Freedom's Dark Horizon, Freedom House. Available at https://freedomhouse.org/report/freedom-press/freedom-press-2017?gclid=CPeczoyF9dQCFUKBswodmrAl2w (Last accessed 6/2/2017)

<sup>37</sup> Alliance for Affordable Internet (2017), Affordability Report. Web Foundation. Available at http://1e8q3q16vyc81g8l3h3md6q5f5e.wpengine.netdna-cdn.com/wp-content/ uploads/2017/02/A4Al-2017-Affordability-Report.pdf (Last accessed 5/2/2017)

Graham, M. (2014, October 29). Dominant Wikipedia language by country. Mark Graham's Blog. Available at http://www.markgraham.space/blog/dominant-wikipedia-language-by-country (Last accessed 5/2/2017); and Noack, R. & Gamio, L (2015, April 23). The world's languages, in 7 maps and charts. Washington Post. Available at https://www.washingtonpost.com/news/worldviews/wp/2015/04/23/the-worlds-languages-in-7-maps-and-charts/?utm\_term=.ff59ea679873 (Last accessed 5/6/2017)
 World Bank Databank (2017). World Bank. Available at http://data.worldbank.org/indicator/SE.ADT.LITR.ZS?locations=XO-XD (Last accessed 5/2/2017)

<sup>41</sup> Kafka, P (2017, March 4) Facebook has started to flag fake news stories. recode. Available at https://www.recode.net/2017/3/4/14816254/facebook-fake-news-disputedtrump-snopes-politifact-seattle-tribune (Last accessed 5/2/2017); https://www.theguardian.com/technology/2017/mar/22/facebook-fact-checking-tool-fake-news (Last accessed 5/2/2017)

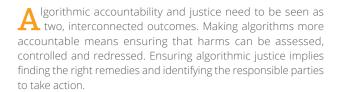
<sup>42</sup> boyd, d. (2017, March 7) Google and Facebook can't just make fake news disappear. Backchannel Available at https://backchannel.com/google-and-facebook-cant-just-make-fake-news-disappear-48f4b4e5fbe8 (Last accessed 5/2/2017)

<sup>43</sup> Sigal, I. (2017, March 20). Fake News and Fake Solutions: How Do We Build a Civics of Trust? Global Voices Advox. Available at https://advox.globalvoices.org/2017/03/20/ fake-news-and-fake-solutions-how-do-we-build-a-civics-of-trust/ (Last accessed 5/2/2017)

<sup>44</sup> Jones, C (2017, March n.d.) What will FG do to people who publish fake news. Naij. Available at https://www.naij.com/1090981-what-fg-people-publish-fake-news-buharismedia-aide-reveals-9-crucial-facts.html (Last accessed 5/2/2017)



## A WAY FORWARD



Achieving accountability in the age of algorithms can be seen as a quality assessment indicator.<sup>45</sup> This can help to provide evidence of the presence of potential harms within the complex systems of algorithms. Based on the different types of harms that have been identified, different actors then can work collectively together to determine what kind of policy, ethical and technical levers need to be pulled to repair these problems.

Achieving algorithmic justice will equally require addressing a range of technical, ethical, policy and knowledge gaps. Implementing these solutions requires a recognition that this is a shared responsibility of all stakeholders: algorithmic system designers, legal and regulatory authorities, public interest groups and users.

Approaches by different actors to achieve algorithmic accountability and justice are nascent and mainly have been focused on the US, UK and Europe. But the suggested actions can also be valid for application in other contexts, including low and middle-income countries (see below). They point to the need to take a systemic approach to the problem, where geographical differences are less important than understanding and accounting for the mechanics behind algorithms. What is still essential, however, is to better document the use and potential harms of algorithms in low and middle-income countries, which to date have been seldom studied.

Society is coming to terms that algorithms impact — in positive and negative ways — our everyday interactions. The question now is how to make sure that we collectively better design them, understand them, remedy any harms, and ensure that the expected positive outcomes of algorithms are realised and evenly distributed.

45 Interview with Anupam Datta (Carnegie Mellon University)

### POTENTIAL AREAS FOR ACTION

### 1. Advocacy groups and public interest organisations:

- Engage citizens and governments on the risks associated with algorithmic decision-making.
- Take a principles-first approach in organising workshops and conferences within sectors to determine the fundamental values that should be embedded in public sector algorithms.
- Act as a bridge between the technology and government sector.
- Advocate for the active involvement of traditionally excluded groups in the process of designing algorithms to ensure values are appropriately translated into code.
- Actively take a stand against instances of unfair and/ or biased assessment and treatment.
- Promote the availability of fair, open data sets for training algorithmic models.
- Organise context-specific algorithmic accountability workshops.

# 2. Companies and other data owners:

- Invest in quality controls to oversee data collection processes.
- Ensure human-in-the-loop verification (e.g. involving human operators within automated decision systems).
- Be transparent in communicating internal processes and accounting options with the public as soon as algorithmic harms are claimed and are being evaluated.
- Participate in stakeholder discussion in preventive efforts, as well as efforts towards repair.
- Define and promote the use of a code of conduct for responsible use of data and algorithms.
- Push towards transforming the entire industry into being more open and accountable for their data and services.
- Intensify the diversity, equity and inclusivity (DEI) efforts to go beyond human resources, and allowed to effectively influence the approach towards product development, and services provision by the company.

### 3. Foundations:

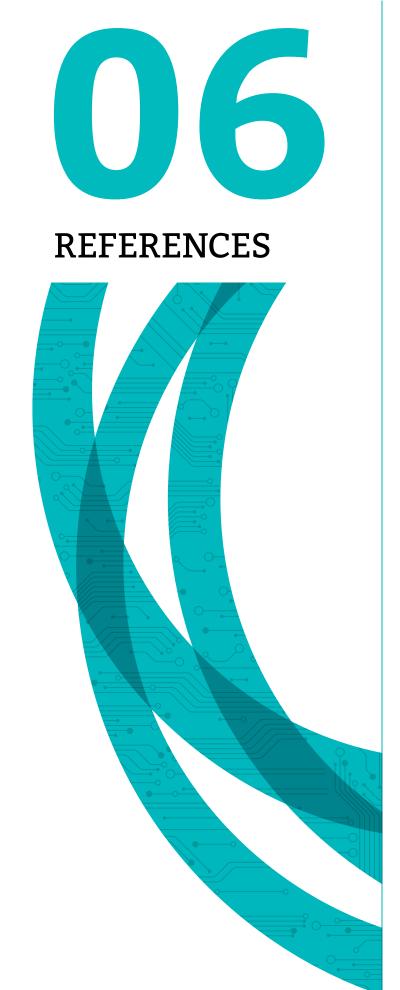
- Invest in technical research to promote more viable options for accounting.
- Coordinate multi-stakeholder engagements and efforts in repair options requiring international and cross-sector cooperation. Consider focusing on specific algorithm-related industries and sectors for greater context and impact.
- Provide platforms for researchers to help craft toolkits and evaluation frameworks for governments and companies purchasing third party algorithms or data services.
- Promote models and solutions that facilitate inclusion of underserved, underrepresented groups.

# 4. Policymakers in national governments:

- Promote government research in computer science and legal innovations.
- Establish active platforms to engage decision-makers in balanced debates regarding the opportunities and risks posed algorithms to society, including legislative research and inquiries.
- Promote the definition and verification of standards.
- Pursue sector regulation and policy reform.

### 5. Universities:

- Coordinate cross-departmental academic workshops
  and meetings.
- Promote technical research, experimentation, and collaboration among academics, companies and governments.
- Ensure ethics courses are part of the basic computer science curricula, and that basic knowledge of code is included in the social sciences curricula.
- Facilitate building models and platforms of trust and collaboration.



- Angwin, J., Larson, J., Mattu, S., Kirchner, L.: Machine bias. ProPublica (2016). URL https://www.propublica.org/article/ machine-bias-risk-assessments-in- criminal-sentencing Barocas, S., Selbst, A.: Big data's disparate impact. California Law Review 104, 671–732 (2016)
- Barry-Jester, A.M., Casselman, B., Goldstein, D.: The new science of sentencing. The Marshall Project (2015). URL https://www.themarshallproject.org/2015/08/04/the-newscience-of-sentencing.
- Burrell, J.: How the machine 'thinks': Understanding opacity in machine learning algorithms. Big Data & Society 3(1) (2016)
- Calders, T., Verwer, S.: Three naive bayes approaches for discrimination-free classification. Data Mining and Knowledge Discovery 21(2), 277–292 (2010)
- Calders, T., Zliobaite, I.: Why unbiased computational processes can lead to discriminative decision procedures.
   In: B. Custers, T. Calders, B. Schermer, T. Zarsky (eds.)
   Discrimination and Privacy in the Information Society, pp. 43–57 (2013)
- Chouldechova, S.: Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. arXiv preprint arXiv:1610.07524 (2016)
- Christin, A., Rosenblatt, A., boyd, d.: Courts and predictive algorithms. Data & Civil Rights Primer (2015)
- Citron, D., Pasquale, F.: The scored society. Washington Law Review 89(1), 1–33 (2014)
- Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., Huq, A.: Fair algorithms and the equal treatment principle. Working Paper (2017)
- Crawford, K., Schultz, J.: Big data and due process: Toward a framework to redress predictive privacy harms. Boston College Law Review 55(1), 93–128 (2014)
- Datta, A., Tschantz, M.C.: Automated experiments on ad privacy settings. In: Proceedings on Privacy Enhancing Technologies, pp. 92–112 (2015)
- Diakopoulos, N. Algorithmic accountability: Journalistic investigation of computational power structures. Digital Journalism (2015)
- D'Ignazio, C. & Bhargava, R.: Approaches to Building Big Data Literacy. MIT Media Lab (2015) https://dam-prod. media.mit.edu/x/2016/10/20/Edu\_D'Ignazio\_52.pdf Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.: Fairness through awareness. In: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, pp. 214–226. ACM (2012)
- Feldman, M., Friedler, S., Moeller, J., Scheidegger, C., Venkatasubramanian, S.: Certifying and removing disparate impact. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 259–268 (2015)
- Fiske, S.: Stereotyping, prejudice, and discrimination. In:
  D. Gilbert, S. Fiske, G. Lindzey (eds.) Handbook of Social Psychology, pp. 357–411. Boston: McGraw-Hill (1998)

- Friedler, S.A., Scheidegger, C., Venkatasubramanian,
  S.: On the (im)possibility of fairness. arXiv preprint arXiv:1609.07236 (2016)
- Gillespie, T.: The relevance of algorithms. In: T. Gillespie, P. Boczkowski, K. Foot (eds.) Media technologies: Essays on communication, materiality, and society, pp. 167–193. MIT Press (2014)
- Hajian, S., Bonchi, F., Castillo, C.: Algorithmic bias: From discrimination discovery to fairness-aware data mining. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 2125–2126. ACM (2016)
- Kamiran, F., Calders, T., Pechenizkiy, M.: Discrimination aware decision tree learning. In: Proceedings of 2010 IEEE International Conference on Data Mining, pp. 869–874. IEEE (2010)
- Kleinberg, J., Mullainathan, S., Raghavan, M.: Inherent tradeoffs in the fair determination of risk scores. In: Proceedings of the 8th Innovations in Theoretical Computer Science Conference. ACM (2017)
- Kroll, J.A., Huey, J., Barocas, S., Felten, E.W., Reidenberg, J.R., Robinson, D.G., Yu, H.: Accountable algorithms. University of Pennsylvania Law Review 165 (2017)
- Lepri, Bruno, Staiano, J., Sangokoya, D., Letouzé, E., & Oliver, N. (2017). "The Tyranny of Data? The Bright and Dark Sides of Data-Driven Decision-Making for Social Good." Transparent Data Mining for Big and Small Data. Springer International Publishing. Macnish, K.: Unblinking eyes: The ethics of automating surveillance. Ethics and Information Technology 14(2), 151–167 (2012)
- O'Neil, C.: Weapons of math destruction: How big data increases inequality and threatens democracy. Crown (2016)
- Pager, D., Shepherd, H.: The sociology of discrimination: Racial discrimination in employment, housing, credit and consumer market. Annual Review of Sociology 34, 181–209 (2008)
- Pasquale, F.: The Black Box Society: The secret algorithms that control money and information. Harvard University Press (2015)
- Podesta, J., Pritzker, P., Moniz, E., Holdren, J., Zients, J.: Big data: Seizing opportunities, preserving values. Tech. rep., Executive Office of the President (2014)
- Ramirez, E., Brill, J., Ohlhausen, M., McSweeney, T.: Big data: A tool for inclusion or exclusion? Tech. rep., Federal Trade Commission (2016)
- Sandvig, C., Hamilton, K., Karahalios, K., Langbort, C. : Auditing algorithms: Research methods for detecting discrimination on internet platforms. In: Data and Discrimination: Converting Critical Concerns into Productive Inquiry, a preconference at the 64th Annual Meeting of the International Communication Association (2014)
- Schermer, B.W.: The limits of privacy in automated profiling and data mining. Computer Law & Security Review 27(1), 45–52 (2011)

- Tobler, C.: Limits and potential of the concept of indirect discrimination. Tech. rep., European Network of Legal Experts in Anti-Discrimination (2008)
- Wang, T., Rudin, C., Wagner, D., Sevieri, R.: Learning to detect patterns of crime. In: Machine Learning and Knowledge Discovery in Databases, pp. 515–530. Springer (2013)
- Willson, M.: Algorithms (and the) everyday. Information, Communication & Society (2016)
- Zafar, M.B., Martinez, I.V., Rodriguez, M.D., Gummadi, K.P.: Learning fair classifiers. arXiv preprint arXiv:1507.05259 (2015)
- Zarsky, T.: The trouble with algorithmic decisions: An analytic road map to examine efficiency and fairness in automated and opaque decision making. Science, Technology, and Human Values 41(1), 118–132 (2016)
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., Dwork, C.: Learning fair representation. In: Proceedings of the 2013 International Conference on Machine Learning (ICML), pp. 325–333 (2012)

Interviews for this study were conducted by David Sangokoya, with assistance from Gabriel Pestre (Data-Pop Alliance), with the following individuals: Robyn Caplan (Data and Society Research Institute), Anupam Datta (Carnegie Mellon University), Sorelle Friedler (Haverford University), Yves-Alexandre de Montjoye (Imperial University), Daniel Neill (Carnegie Mellon University / New York University), Arnaud Sahuguet (Cornell Tech), Ravi Shroff (New York University) and Suresh Venkatasubramanian (University of Utah).



World Wide Web Foundation, 1110 Vermont Ave NW, Suite 500, Washington DC 20005, USA

www.webfoundation.org | Twitter: @webfoundation