

# GigaScience

## Where is the human in the data? A guide to ethical data use

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-18-00208	
<b>Full Title:</b>	Where is the human in the data? A guide to ethical data use	
<b>Article Type:</b>	Commentary	
<b>Funding Information:</b>	Marsden Fund Fast Start Grant (UOO1515)	Dr Angela Ballantyne
<b>Abstract:</b>	Being asked to write about the ethics of big data is a bit like being asked to write about the ethics of life. Big data is now integral to so many aspects of our daily lives – communication, social interaction, medicine, access to government services, shopping and navigation. Given this diversity, there is no one-size-fits-all framework for how to ethically manage your data.	
<b>Corresponding Author:</b>	Angela Ballantyne  NEW ZEALAND	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Angela Ballantyne	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Angela Ballantyne	
<b>Order of Authors Secondary Information:</b>		
<b>Additional Information:</b>		
<b>Question</b>	<b>Response</b>	
Are you submitting this manuscript to a special series or article collection?	No	
<b>Experimental design and statistics</b>	Yes	
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.		
Have you included all the information requested in your manuscript?		
<b>Resources</b>	Yes	
A description of all resources used, including antibodies, cell lines, animals and software tools, with enough		

<p>information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

[Click here to view linked References](#)

FINAL 01 JUNE 2018

## Title: “Where is the human in the data? A guide to ethical data use”

Angela Ballantyne<sup>1,\*</sup>

1. Department of Primary Health Care and General Practice at University of Otago, Wellington, New Zealand.

\*corresponding author

Email: [angela.ballantyne@otago.ac.nz](mailto:angela.ballantyne@otago.ac.nz)

ORCID: 0000-0003-2666-9557

Keywords: Ethics, Big Data, Data Analytics, Data Protection

### Abstract:

Being asked to write about the ethics of big data is a bit like being asked to write about the ethics of life. Big data is now integral to so many aspects of our daily lives – communication, social interaction, medicine, access to government services, shopping and navigation. Given this diversity, there is no one-size-fits-all framework for how to ethically manage your data. With that in mind, here I attempt to present seven ethical values for responsible data use.

### Body Text:

Data is ubiquitous because it is so useful. This means that many different parties – data subjects and sources, associated communities, researchers, governments and businesses – will have competing interests in relation to the data. Just as we make trade-offs in our daily life (to walk or to drive to work? doughnut versus salad for lunch?) we need to make trade-offs about competing interests in relation to data.

I am talking here about interests, rather than rights. Note that many parties who don't have legal rights to control access to and use of data, may none-the-less have compelling interests in the data. Responsible data use requires attention to these broad interests. Facebook's recent troubles highlight this. Even if Facebook was legally entitled to share users' data with Cambridge Analytica, Facebook massively under-estimated users' interests and expectations in relation to privacy, control and appropriate use.

In areas of rapid progress, such as data science, practice can quickly outstrip the legal framework. Data use may be within the parameters of the law (e.g. data protection or privacy regulation), but may nonetheless be unethical and/or outside the social licence.<sup>i</sup> We should be aiming to align the social licence, ethics and the law to ensure that data use is publicly acceptable, normatively justified and legal. Where there is misalignment of the law, ethics, and the social licence, data users need to tread carefully.

The following is a list of ethical values, also depicted in Figure 1, that can: (1) help identify who has an interest in the data and where these interests might clash; (2) help data holders to articulate the

---

<sup>i</sup> 'Ethics' is normative – it makes a claim about what the morally correct course of action would be, and attributes praise or blame. 'Social licence' is descriptive – it describes whether a given data use is accepted by the data subjects, public, and other stakeholders.

1 ethical trade-offs that need to be made; and (3) guide deliberation about responsible data use. The  
2 values often clash – maximising data security will conflict with maximising social value through  
3 broader data use. In different circumstances priority will appropriately be given to different values.  
4 This process is about making informed, explicit and justifiable trade-offs, rather than following a set  
5 of prescribed rules.  
6

7 **Social value** Data is in demand because it has value. Data can contribute to knowledge and  
8 innovation, drive efficiency, reduce harm from ineffective or poorly targeted services and reduce  
9 costs. Open data is important to drive the advancement of scientific knowledge, preserve datasets,  
10 test and verify conclusions, refine algorithms, and safeguard against misconduct.  
11

12 **Harm minimization** Data collection, storage and use should be designed to minimise and manage  
13 risks of harm. Harms can be physical, economic, psychological or reputational and can be  
14 experienced by individuals, communities or organisations. Anonymization (pseudonymization and  
15 de-identification) has been the cornerstone of protecting individual data subjects from harm. But  
16 anonymization is failing in the era of big data, where there are hundreds of thousands of data points  
17 for a single individual.<sup>1</sup> Data scientists have proved repeatedly that they can re-identify individuals in  
18 supposedly anonymous data sets.<sup>2</sup> Furthermore, anonymization and de-identification do little to  
19 protect communities from harm. Data analytics and AI are increasingly used to characterize the  
20 behaviour of communities and inform the delivery of services. Data can be used to stigmatize or  
21 discriminate.  
22

23 **Control** Control refers to the capacity for data-subjects to be autonomous and self-determining.  
24 Were data subjects asked for their consent at the point of data collection? To what degree will data  
25 subjects' preferences determine how the data is used? Is this a secondary use of the data that differs  
26 from the original consent? Is the data use novel and original or is it likely to be consistent with the  
27 expectations' of data subjects? Various models of consent have been proposed for data – including  
28 broad consent, dynamic consent<sup>3</sup> and meta-consent<sup>4</sup>. However much data use (especially linking and  
29 secondary uses) occurs without consent. In these cases, data users need to be safe stewards of the  
30 data. Transparency, engagement and accountability are especially important for data used without  
31 consent.<sup>5</sup>  
32

33 **Justice** Justice concerns the equitable treatment of those with an interest in the data activities;  
34 including the fair distribution of any benefits and burdens arising from the collection, storage, use,  
35 linkage and sharing of data. The term 'benefit sharing' was first used in relation to non-human  
36 genetic resources in the Convention on Biological Diversity (CBD) adopted at the Earth Summit in Rio  
37 de Janeiro, Brazil in 1992. Benefit sharing requires that the advantages/profits derived from the  
38 data are shared fairly amongst the data providers and to the community from which the data  
39 originate. Recent data advocacy, especially in relation to indigenous data, has moved away from  
40 'benefit sharing' towards 'power sharing' – arguing that data subjects and communities should have  
41 decision making capacity in relation to data governance and use.<sup>6</sup>  
42

43 **Trustworthiness** Trustworthiness is the property of being worthy of trust - it can apply to individuals,  
44 organisations, and institutions, but also relates to data quality, systems of knowledge production,  
45 scientific integrity and professional standards.<sup>7</sup> When judging trustworthiness, we look for  
46 truthfulness, reliability, and consistency but also goodwill. A robust data ecosystem requires a high-  
47 level of trust. A breach of trust can affect not only the agents involved, but an entire profession or  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 institution. The dispute between Arizona State University and members of the Havasupai Indian  
2 tribe, over the use of genetic samples for research left a legacy of mistrust and fear of exploitation.<sup>8</sup>  
3 As Smith famously argued “‘Research’ is probably one of the dirtiest words in the indigenous world’s  
4 vocabulary.”<sup>9</sup> And trust, when lost, can take significant efforts to rebuild.<sup>10</sup>  
5

6 **Transparency** Transparency is openness and accessibility in decision-making and actions. When the  
7 data activity occurs without the data subjects’ consent and is justified on the grounds of ‘social  
8 value’, the arguments in favour of transparency and openness are especially compelling.  
9

10 Transparency helps to demonstrate respect for data subjects and trustworthiness; and it underpins  
11 public engagement and accountability. Full transparency would include a public description of the  
12 data activity, purpose and justification, anticipated social value, harm-mitigation strategies, public  
13 engagement strategies, level of security and encryption, research results, and the coding/algorithms.  
14 When launching a £1.5 Billion initiative in AI in April 2018, President Macron announced that anyone  
15 receiving AI funding money from the French government will be required to make their algorithms  
16 open and transparent.  
17  
18  
19

20 **Accountability** Accountability refers to holding data users and custodians responsible for the  
21 consequences of their decisions and actions. Data regulation is increasingly focused on  
22 accountability. A significant innovation in the EU General Data Protection Regulation (GDPR) (which  
23 will come into force on the 25th May 2018) is the introduction of ‘accountability’ (Article 5(2)) to the  
24 list of principles relating to personal data. Under the GDPR organisations will need to be more  
25 intentional about their data collection and use and maintain open lines of communication with data  
26 subjects.  
27  
28  
29  
30

31 Given these competing values, there will be multiple different ‘ethical’ solutions to data  
32 management. The task is to identify the ethical issues, reason though how to balance conflicting  
33 demands, articulate the trade-offs and justify the conclusions. Do this as publically and  
34 transparently as possible; and make time to revise and re-assess.  
35  
36

37 We use data to tell stories, to make sense of the world. This means telling stories about *people* and  
38 how they live. Data has the appealing veneer of scientific objectivity; but the process of telling  
39 stories is never ethically neutral. Our starting point should be to ask: Where is the human in the  
40 data? What would this data use look like from the data subjects’ perspective?  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2 **Abbreviations**

3  
4 AI: Artificial Intelligence; CBD: Convention on Biological Diversity; GDPR: General Data Protection  
5 Regulation  
6

7  
8 **Declarations**

9  
10 **Competing interests** The author declares that they have no competing interests.

11  
12 **Funding** Marsden Fund Fast Start Grant (UOO1515) 2016-2019 “The ethics of research on clinical  
13 data and tissue without explicit patient consent”  
14

15  
16 **Acknowledgements**

17  
18 I would like to thank the *Working Group for Big Data Ethics in Health and Research* who are a delight  
19 to work with and who contributed to the definitions of the values used here. The Working Group is  
20 organized by the Science, Health and Policy-Relevant Ethics (SHAPES) initiative at the Centre for  
21 Biomedical Ethics, Yong Loo Lin School of Medicine at the National University of Singapore. Co-  
22 Chairs: Prof Tai E Shyong (National University of Singapore/National University Hospital) and Prof  
23 Graeme Laurie (University of Edinburgh). Members: Dr Angela Ballantyne (University of Otago); Dr  
24 Iain Brassington (University of Manchester); Mr Markus Labude (National University of Singapore);  
25 A/Prof Hannah Lim Yee Fen (Nanyang Technological University); A/Prof Wendy Lipworth (University  
26 of Sydney); Dr Tamra Lysaght (National University of Singapore); Dr Owen Schaefer (National  
27 University of Singapore); Prof Cameron Stewart (University of Sydney); A/Prof Shirley Sun Hsiao-Li  
28 (Nanyang Technological University); and Dr Vicki Xafis (National University of Singapore).  
29  
30  
31  
32  
33

34 Thanks also to those who attended the Data Ethics Governance Workshop in Wellington NZ in 2017,  
35 for fruitful discussion of these topics, especially Maria Stubbe, June Atkinson and Rochelle Style.  
36

37  
38 **Figure legend**

39  
40 **Figure 1: An infographic summarizing the ethical values**  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**References.**

- 1  
2  
3 <sup>1</sup> Roy M. Data anonymization techniques less reliable in era of big data. TechTarget. 2017.  
4 [https://searchcompliance.techtarget.com/feature/High-dimensional-info-complicates-data-anonymization-](https://searchcompliance.techtarget.com/feature/High-dimensional-info-complicates-data-anonymization-techniques)  
5 [techniques](https://searchcompliance.techtarget.com/feature/High-dimensional-info-complicates-data-anonymization-techniques) Accessed 10 May 2018.
- 6 <sup>2</sup> Montjoye AJ, Radaelli L, Singh VK, Pentland A. Unique in the shopping mall: On the reidentifiability of credit  
7 card metadata. *Science*. 2015; 536-539.
- 8 <sup>3</sup> Kaye J, Whitley EA, Lund D, Morrison M, Teare H, Melham K. 2015. Dynamic consent: a patient interface for  
9 twenty-first century research networks. *Eur J Hum Genet*. 2015 Feb;23(2):141-6. doi: 10.1038/ejhg.2014.71
- 10 <sup>4</sup> Ploug T, Holm S.. Meta Consent - A Flexible Solution to the Problem of Secondary Use of Health Data.  
11 *Bioethics*. 2016; 30(9):721-732. doi: 10.1111/bioe.12286.
- 12 <sup>5</sup> Ballantyne A, Schaefer GO. Consent and the ethical duty to participate in health data research. *J Med Ethics*.  
13 2018. doi: 10.1136/medethics-2017-104550.
- 14 <sup>6</sup> Kukutai T, Taylor J. Data sovereignty for indigenous peoples: current practice and future needs. In T. Kukutai,  
15 & J. Taylor (Eds.), *Indigenous Data Sovereignty: toward an agenda* (pp. 1-22). 2016. Acton, Australia: ANU  
16 press. Retrieved from <https://press.anu.edu.au/>
- 17 <sup>7</sup> Aitken M, Cunningham-Burley S, Pagliari C. Moving from trust to trustworthiness: Experiences of public  
18 engagement in the Scottish Health Informatics Programme. *Science & Public Policy*. 2016; 43(5):713-723.  
19 doi:10.1093/scipol/scv075.
- 20 <sup>8</sup> Mello MM, Wolf LE. The Havasupai Indian tribe case--lessons for research involving stored biologic samples.  
21 *N Engl J Med*. 2010 Jul 15; 363(3):204-7. doi: 10.1056/NEJMp1005203.
- 22 <sup>9</sup> Tuhiwai Smith, Linda, 1999. *Decolonizing Methodologies : Research and Indigenous Peoples*. London:  
23 University of Otago Press.
- 24 <sup>10</sup> Carter P, Laurie GT, Dixon-Woods M. The social licence for research: why care.data ran into trouble.  
25 *J Med Ethics*. 2015; 41(5):404-9.
- 26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

