



King's Research Portal

Document Version
Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Brandão, M., Mansouri, M., & Magnusson, M. (2022). Editorial: Responsible Robotics. *Frontiers in Robotics and AI*, 9.

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

Editorial: Responsible Robotics

Martim Brandão^{1,*}, Masoumeh Mansouri² and Martin Magnusson³

¹*Department of Informatics, King's College London, London, United Kingdom*

²*School of Computer Science, University of Birmingham, Birmingham, United Kingdom*

³*School of Science and Technology (AASS), Örebro University, Örebro, Sweden*

Correspondence*:

Martim Brandão

martim.brandao@kcl.ac.uk

2 **Keywords:** Robotics, Responsible Innovation, Responsible Robotics, Trustworthy Robotics, Critical Robotics, AI and Society, Robot
3 **Ethics**

RESPONSIBLE AI AND ROBOTICS

4 Recent work in both academia, industry, and journalism has brought widespread attention to various
5 kinds of harmful impact that AI can have on society. These are very often concentrated on marginalized
6 social groups. AI algorithms may unintentionally reinforce social prejudice Bolukbasi et al. (2016) and
7 biased conceptions of gender Adams and Loideáin (2019); Hamidi et al. (2018), race Sweeney (2013), age
8 Rosales and Fernández-Ardèvol (2019) or disabilities Guo et al. (2020), they may lead to unfair access
9 to opportunities Dastin (2018); Angwin et al. (2016), discriminatory pricing practices Bar-Gill (2019);
10 Hannak et al. (2014), etc. Recent work has also shown that many seemingly technical issues in machine
11 learning are actually socio-technical. For example: the over-fitting of machine learning models, the choice
12 of dataset or learning objective, and other aspects of learning may lead to algorithms performing poorly
13 on unrepresented or unmodeled groups of people Brandao (2019); Barocas et al. (2019); Buolamwini
14 and Gebu (2018). A growing community of Fairness, Accountability, Transparency, and Ethics of AI¹ is
15 now approaching these topics from a socio-technical point-of-view, in order to identify, understand, and
16 alleviate such issues.

17 Robotics, as a technology focused on automation and intelligent behavior, also abounds in similar ethical
18 and social issues that need to be identified, characterized, and considered in design. While many of the
19 same problems with AI will also be present in robotics, the physical nature of robotics raises new aspects
20 of the social and ethical nature of these technologies. As one example: models that are considerably less
21 accurate on certain groups of people can lead to physical safety differentials Brandao (2019), where robots
22 or autonomous vehicles using those models are more likely to collide with those groups. Additionally, there
23 are physical safety concerns with respect to surgical and other medical robots Yang et al. (2017); Ficuciello
24 et al. (2019), as well as concerns of physical and political security—not least concerning autonomous
25 weapon systems and the dual-use of robot technologies like autonomous cars and drones Brundage et al.
26 (2018); Sparrow (2007).

27 The physical design and visual appearance of robots also introduce new aspects to responsible
28 development. For example, people's moral evaluation of robot decisions can be affected by whether

¹ Example venues: ACM Conference on Fairness, Accountability, and Transparency (FAccT), AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES)

29 the robot is more or less human-like Malle et al. (2016), the design of robots in a care setting affects
30 caregivers and caretakers van Wynsberghe (2021); Kubota et al. (2021), the choice of sensors, measurements
31 and motion has an impact of privacy Calo (2011); Eick and Antón (2020); Luo et al. (2020), and the ethics
32 of deception takes on new shape Danaher (2020).

33 The robotics community has been discussing ethics for long². Recent workshops have also started
34 bringing attention to philosophical problems in robotics³ and issues such as bias⁴ and transparency⁵.
35 These efforts share a common goal of developing robotics technologies responsibly—they are part of
36 “Responsible Robotics” or “Trustworthy Robotics”.

37 A similar effort on “Critical Robotics” Serholt et al. (2021) has focused on questioning current practices
38 in robotics research. These range from how older adults are represented in HRI Burema (2021) and ethical
39 issues in education robots Serholt et al. (2017), to normative dimensions of speech used by researchers
40 Brandao (2021), their technological optimism Šabanović (2010) and the influence of their social background
41 in research directions Forsythe (2001); Šabanović (2010).

THIS RESEARCH TOPIC

42 This research topic gathers a diverse set of articles on Responsible Robotics. They range from user
43 studies and philosophical inquiry, to modeling, algorithmic, and governance methods. Our goal when
44 organizing this research topic was exactly to join various approaches in a single edition—to allow for
45 greater multidisciplinary exchange under the common mission of Responsible Robotics. We believe that
46 Responsible Robotics should focus both on *identifying* social and ethical issues, and on *designing* methods
47 to account for (and alleviate) such issues—thus the focus of this edition on both understanding and *acting*
48 on social and ethical issues.

49 Two articles in the research topic are focused on eliciting social and ethical issues *from users and*
50 *stakeholders*. Lutz and Tamò-Larrioux (2021) investigate privacy concerns of lay users and their impact on
51 technology use intentions, when using social robots that are either privacy-friendly or privacy-invasive (e.g.
52 listen to conversations, share data with third parties). Colombino et al. (2021) use ethnographic studies,
53 interviews and futuristic autobiographies to identify organizational principles, potential roles, and ethical
54 design considerations for a robot that collaborates with disabled employees.

55 Three articles are more focused on methods, or socio-technical solutions to ethical problems in robotics.
56 Webb et al. (2021), for example, focus on methods for conducting investigations of accidents involving
57 humans and robots. In particular, they propose and preliminarily evaluate a role-play-based methodology
58 for investigating accidents, and to evaluate the testimonies that humans can give in forensic investigations
59 of such accidents. Hurtado et al. (2021) focus on issues of harmful social bias in robot learning and how
60 they could be detected and alleviated. Namely, they show through various examples how social robot
61 navigation techniques that mimic human behavior may lead to harmful behavior, such as higher intrusion
62 of personal space or longer waiting times for some groups compared to others. Winfield et al. (2021)
63 focus on issues of transparency from a governance perspective. They describe a new draft standard on
64 transparency for autonomous systems, with several contributions such as transparency levels, measurability,
65 stakeholders, and example-based guidance on using the draft standard.

² ICRA 2007/2009/2011 workshops on Roboethics, ICRA 2014 workshop on “Robotics and Military Applications”

³ Robophilosophy Conference

⁴ ICRA 2019 workshops on “Bias-sensitizing robot behaviours” and “Unlearning biases in robot design”

⁵ HRI 2022 workshop on “Fairness and Transparency in HRI”, ICRA 2020 workshop “Against robot dystopias”

66 We then dive into philosophical inquiry and frameworks for robot ethics. Rhim et al. (2021) combine work
67 in moral philosophy and psychology to propose a model that explains human decision-making in moral
68 dilemmas involving autonomous vehicles. Pirni et al. (2021) consider aspects of autonomy and vulnerability
69 in the ethics of designing care robots. And Kuipers (2022) argues that AI and robotics technologies rely
70 heavily on over-simplified models, and that the widespread use of such models can lead to the erosion of
71 trust and cooperation effectiveness. The article can serve as an argument for why more attention should be
72 given to the *modeling* of complex socio-technical factors in AI/robotics.

73 Finally, two articles in the research topic dive into issues of jobs and economics in robotics and automation.
74 Studley (2021) argues that we should consider how robotics impacts global supply chains, international
75 development, and global economic disparities. Kyvik Nordås and Klügl (2021) then use modeling to
76 understand the uptake of automation technologies and its relationship with unemployment and engineering,
77 consultancy, and manufacturing jobs. The authors use this analysis to suggest an automation policy focus
78 on user costs and education.

79 We believe that the contributions collected in this special issue can be relevant to roboticists, AI
80 practitioners, policy makers and any other stakeholders concerned with the societal impacts of AI and
81 robotics. We hope this special issue will stimulate future work on responsible robotics.

82 We end with an important remark. While the abundance of social and ethical issues raised in this
83 editorial and this Research Topic might feel overwhelming or hopeless, we believe the opposite is the case.
84 Responsible Robotics is about clearly identifying potential issues, because by doing so it is also possible
85 to work towards responsible methods that mitigate them. This ultimately facilitates the application of
86 robotics and AI in ways that increase safety, efficiency, and wellbeing in many areas of life: transportation,
87 healthcare, work life, just to name a few.

AUTHOR CONTRIBUTIONS

88 All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it
89 for publication.

REFERENCES

- 90 Adams, R. and Loideáin, N. N. (2019). Addressing indirect discrimination and gender stereotypes in ai
91 virtual personal assistants: the role of international human rights law. *Cambridge International Law*
92 *Journal* 8, 241–257
- 93 Angwin, J., Larson, J., Mattu, S., and Kirchner, L. (2016). Machine bias: there's software used across the
94 country to predict future criminals. and it's biased against blacks.
- 95 Bar-Gill, O. (2019). Algorithmic price discrimination when demand is a function of both preferences and
96 (mis) perceptions. *University of Chicago Law Review* 86
- 97 Barocas, S., Hardt, M., and Narayanan, A. (2019). *Fairness and Machine Learning* (fairmlbook.org).
98 <http://www.fairmlbook.org>
- 99 Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016). Man is to computer
100 programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information*
101 *processing systems* 29
- 102 Brandao, M. (2019). Age and gender bias in pedestrian detection algorithms. In *Workshop on Fairness*
103 *Accountability Transparency and Ethics in Computer Vision, CVPR*

- 104 Brandao, M. (2021). Normative roboticists: the visions and values of technical robotics papers. In
105 *IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*. 671–677.
106 doi:10.1109/RO-MAN50785.2021.9515504
- 107 Brundage, M., Avin, S., Clark, J., Toner, H., Eckersley, P., Garfinkel, B., et al. (2018). The malicious use
108 of artificial intelligence: Forecasting, prevention, and mitigation. *arXiv preprint arXiv:1802.07228*
- 109 Buolamwini, J. and Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial
110 gender classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*,
111 eds. S. A. Friedler and C. Wilson (New York, NY, USA: PMLR), vol. 81 of *Proceedings of Machine*
112 *Learning Research*, 77–91
- 113 Burema, D. (2021). A critical analysis of the representations of older adults in the field of human–robot
114 interaction. *AI & Society* , 1–11
- 115 Calo, R. (2011). Robots and privacy. *Robot Ethics: The Ethical and Social Implications of Robotics*
- 116 Colombino, T., Gallo, D., Shreepriya, S., Im, Y., and Cha, S. (2021). Ethical design of a robot platform
117 for disabled employees: Some practical methodological considerations. *Frontiers in Robotics and AI* 8.
118 doi:10.3389/frobt.2021.643160
- 119 Danaher, J. (2020). Robot betrayal: a guide to the ethics of robotic deception. *Ethics and Information*
120 *Technology* 22, 117–128
- 121 Dastin, J. (2018). Amazon scraps secret ai recruiting tool that showed bias against women. In *Ethics of*
122 *Data and Analytics* (Auerbach Publications). 296–299
- 123 Eick, S. and Antón, A. I. (2020). Enhancing privacy in robotics via judicious sensor selection. In *2020 IEEE*
124 *International Conference on Robotics and Automation (ICRA)*. 7156–7165. doi:10.1109/ICRA40945.
125 2020.9196983
- 126 Ficuciello, F., Tamburrini, G., Arezzo, A., Villani, L., and Siciliano, B. (2019). Autonomy in surgical
127 robots and its meaningful human control. *Paladyn, Journal of Behavioral Robotics* 10, 30–43. doi:doi:
128 10.1515/pjbr-2019-0002
- 129 Forsythe, D. (2001). *Studying those who study us: An anthropologist in the world of artificial intelligence*
130 (Stanford University Press)
- 131 Guo, A., Kamar, E., Vaughan, J. W., Wallach, H., and Morris, M. R. (2020). Toward fairness in ai for
132 people with disabilities: a research roadmap. *ACM SIGACCESS Accessibility and Computing* , 1–1
- 133 Hamidi, F., Scheuerman, M. K., and Branham, S. M. (2018). Gender recognition or gender reductionism?
134 the social implications of embedded gender recognition systems. In *Proceedings of the 2018 chi*
135 *conference on human factors in computing systems*. 1–13
- 136 Hannak, A., Soeller, G., Lazer, D., Mislove, A., and Wilson, C. (2014). Measuring price discrimination
137 and steering on e-commerce web sites. In *Proceedings of the 2014 conference on internet measurement*
138 *conference*. 305–318
- 139 Hurtado, J. V., Londoño, L., and Valada, A. (2021). From learning to relearning: A framework for
140 diminishing bias in social robot navigation. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.
141 650325
- 142 Kubota, A., Pourebadi, M., Banh, S., Kim, S., and Riek, L. (2021). Somebody that i used to know: The
143 risks of personalizing robots for dementia care. *Proceedings of We Robot*
- 144 Kuipers, B. (2022). Trust and cooperation. *Frontiers in Robotics and AI* 9. doi:10.3389/frobt.2022.676767
- 145 Kyvik Nordås, H. and Klügl, F. (2021). Drivers of automation and consequences for jobs in engineering
146 services: An agent-based modelling approach. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.
147 637125

- 148 Luo, Y., Yu, Y., Jin, Z., Li, Y., Ding, Z., Zhou, Y., et al. (2020). Privacy-aware uav flights through
149 self-configuring motion planning. In *2020 IEEE International Conference on Robotics and Automation*
150 (*ICRA*). 1169–1175. doi:10.1109/ICRA40945.2020.9197564
- 151 Lutz, C. and Tamò-Larrieux, A. (2021). Do privacy concerns about social robots affect use intentions?
152 evidence from an experimental vignette study. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.
153 627958
- 154 Malle, B. F., Scheutz, M., Forlizzi, J., and Voiklis, J. (2016). Which robot am i thinking about? the
155 impact of action and appearance on people's evaluations of a moral robot. In *2016 11th ACM/IEEE*
156 *International Conference on Human-Robot Interaction (HRI)* (IEEE), 125–132
- 157 Pirni, A., Balistreri, M., Capasso, M., Umbrello, S., and Merenda, F. (2021). Robot care ethics between
158 autonomy and vulnerability: Coupling principles and practices in autonomous systems for care. *Frontiers*
159 *in Robotics and AI* 8. doi:10.3389/frobt.2021.654298
- 160 Rhim, J., Lee, J.-H., Chen, M., and Lim, A. (2021). A deeper look at autonomous vehicle ethics: An
161 integrative ethical decision-making framework to explain moral pluralism. *Frontiers in Robotics and AI*
162 8. doi:10.3389/frobt.2021.632394
- 163 Rosales, A. and Fernández-Ardèvol, M. (2019). Structural ageism in big data approaches. *Nordicom*
164 *Review* 40, 51–64
- 165 Šabanović, S. (2010). Robots in society, society in robots. *International Journal of Social Robotics* 2,
166 439–450
- 167 Serholt, S., Barendregt, W., Vasalou, A., Alves-Oliveira, P., Jones, A., Petisca, S., et al. (2017). The case of
168 classroom robots: teachers' deliberations on the ethical tensions. *AI & Society* 32, 613–631
- 169 Serholt, S., Ljungblad, S., and Ní Bhroin, N. (2021). Introduction: special issue—critical robotics research.
170 *AI & Society* , 1–7
- 171 Sparrow, R. (2007). Killer robots. *Journal of applied philosophy* 24, 62–77
- 172 Studley, M. (2021). Onshoring through automation; perpetuating inequality? *Frontiers in Robotics and AI*
173 8. doi:10.3389/frobt.2021.634297
- 174 Sweeney, L. (2013). Discrimination in online ad delivery. *Communications of the ACM* 56, 44–54
- 175 van Wynsberghe, A. (2021). Social robots and the risks to reciprocity. *AI & Society* , 1–7
- 176 Webb, H., Dumitru, M., van Maris, A., Winkle, K., Jirotko, M., and Winfield, A. (2021). Role-play
177 as responsible robotics: The virtual witness testimony role-play interview for investigating hazardous
178 human-robot interactions. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.644336
- 179 Winfield, A. F. T., Booth, S., Dennis, L. A., Egawa, T., Hastie, H., Jacobs, N., et al. (2021). Ieee p7001: A
180 proposed standard on transparency. *Frontiers in Robotics and AI* 8. doi:10.3389/frobt.2021.665729
- 181 Yang, G.-Z., Cambias, J., Cleary, K., Daimler, E., Drake, J., Dupont, P. E., et al. (2017). Medical
182 robotics—regulatory, ethical, and legal considerations for increasing levels of autonomy. *Science*
183 *Robotics* 2, eaam8638. doi:10.1126/scirobotics.aam8638