

# Human autonomy in the age of artificial intelligence

Current AI policy recommendations differ on what the risks to human autonomy are. To systematically address risks to autonomy, we need to confront the complexity of the concept itself and adapt governance solutions accordingly.

Carina Prunkl

It is hard to overstate the important role that autonomy plays for our moral and political institutions. A cornerstone of human dignity and a prerequisite of liberal democracy, autonomy is often considered a fundamental human value<sup>1–4</sup>. Progress in the development of artificial intelligence (AI) opens up new opportunities for supporting and fostering autonomy, but it simultaneously poses significant risks. Recent incidents of AI-facilitated deception, manipulation or coercion suggest that AI technologies could seriously interfere with human autonomy on a large scale. Cambridge Analytica's attempt to manipulate voters is just one example<sup>5</sup>. Facebook's "emotional contagion" experiment, in which users were swayed towards adopting certain emotional states, is another<sup>6</sup>.

Consequently, human autonomy has become a central theme across guidelines and principles on the responsible development of AI. The European Commission's High-Level Expert Group (HLEG) lists 'respect for autonomy' as the first of its four key ethical principles in its Guidelines on Trustworthy AI<sup>7</sup>. Several other policy documents, including the Association for Computing Machinery's Code of Ethics<sup>8</sup>, the Montreal Declaration for Responsible Development of Artificial Intelligence<sup>9</sup> and the European Commission's White Paper on Artificial Intelligence<sup>10</sup>, equally emphasise the need to protect and respect autonomy, and the Organisation for Economic Co-operation and Development (OECD) lists autonomy as one of its human-centred values<sup>11</sup>.

Despite this frequent call for the protection of autonomy, there remains substantial ambiguity within these documents as to (i) what exactly is meant by the term 'autonomy', as well as (ii) what the risks from AI to autonomy are. In some cases, 'autonomy' remains undefined<sup>8,10</sup>. Often, however, guidelines take different approaches to what they consider the protection of human autonomy to entail. For example, the HLEG advocates that it

entails no "unjustified coercion, deception, or manipulation" by AI systems<sup>7</sup>; the OECD promotes "capacity for human determination"<sup>11</sup>. Others emphasise that "control over and knowledge about autonomous systems"<sup>12</sup> is needed, and yet others stress that principles of human autonomy translate into the protection of "human decision-making power"<sup>13</sup>. This is also consistent with findings by Fjeld et al., who found that autonomy typically provides the theoretical grounding for principles of "human control of technology"<sup>14</sup>.

The result of this heterogeneity is a patchwork of seemingly disjoint policy recommendations. To illustrate this point further: it is one thing to implement measures that protect users from fraudulent online manipulation (e.g., to prevent incidents like the Cambridge Analytica affair), but an entirely different set of measures are required to ensure human decision-making power (e.g., to ensure that the passenger of a driverless car has authority over most of the car's functions). This poses a challenge to policy-makers: How can we adequately address potential risks to human autonomy?

The overall lack of structure in the current discourse threatens to undermine ongoing governance efforts—efforts that are already straining under the complexity of the technical landscape and the large uncertainty of AI's social impacts. Although there has been remarkable scholarly progress in individual areas, such as online manipulation<sup>5,15–18</sup> or healthcare<sup>19,20</sup>, few scholars have discussed the concept of autonomy within a broader technological context<sup>21–23</sup>. To adequately address the risks that AI might pose to human autonomy, we first need a clearer view of *what we mean* by 'human autonomy' and *how* AI technologies could interfere with it. The following aims to add structure to the debate by highlighting different dimensions of human autonomy, providing examples of how AI systems might interfere with them and discussing some of the policy implications

## Human autonomy as agency and authenticity

'Autonomy' is a notoriously complex concept<sup>24,25</sup>, but it generally can be taken to refer to a person's effective capacity for self-governance. This means that the person can act on the basis of beliefs, values, motivations and reasons that are in some relevant sense their own<sup>3,25</sup>. There are (at least) two fundamental aspects to this definition, each pointing to a different set of conditions that need to be fulfilled for a person (or action) to count as autonomous:

1. *Authenticity*. The beliefs, values, motivations and reasons held by a person are in a relevant sense *authentic* to that person, i.e., not the product of external manipulative or distorting influences.
2. *Agency*. A person *is able to act* on the beliefs and values they hold. This implies that they have meaningful options available to them, allowing them to make choices that are of practical import to their life.

Distinguishing between authenticity and agency explains and clarifies some of the heterogeneity found in the current policy discourse. Those calling for protection from AI-facilitated manipulation and deception are primarily addressing the authenticity dimension of autonomy, whereas those emphasising the importance of retaining control over one's own decisions do so in reference to agency.

Here are some explicit examples of how AI systems could affect *authenticity*:

**Manipulation** is a form of external—often covert—influence by which people's decision-making vulnerabilities are targeted and exploited<sup>26</sup>. Through the analysis of large amounts of data, AI systems are able to identify such vulnerabilities and could be used to exploit them. Recommendation systems, often used by search engines and social media platforms, currently pose

one of the highest risks for AI-facilitated online manipulation<sup>5,15–18</sup>.

**Adaptive preference formation** refers to the process of a person adapting their preferences to match the options that are available to them<sup>27</sup>. The increasing use of recommendation algorithms to pre-select online content or options can lead to such adapted preferences, as first studies suggest<sup>28</sup>. This phenomenon might be reinforced by automation bias, the tendency of humans to favour suggestions from computational systems.

**Deception and adaptive belief formation** are another way in which AI systems might affect authenticity, which relies on the availability of adequate information so as to make appropriate judgments. The amplification of conspiratorial content on social media platforms as a result of algorithmic content selection is an example of how AI systems participate in the shaping of beliefs.

Agency, on the other hand, might be negatively affected by the following:

**Loss of opportunities.** AI systems may create new opportunities for individuals to thrive, but they can also lead to a loss of opportunities, such as when automated decision-making algorithms are racially biased and prevent individuals from accessing health care<sup>29</sup>.

**Loss of freedom.** AI might equally contribute to the restriction of basic liberties directly, e.g., through the deployment of military drones, or indirectly, e.g., through the enabling of large-scale surveillance.

**Loss of competence** to make decisions might occur if more and more tasks are routinely outsourced to AI systems, including decision-making in social, medical or financial settings.

**Paternalism** involves well-intentioned infringements on a person's autonomy against their will<sup>30</sup>. AI systems that engage in full paternalistic behaviour are mostly future talk at this point, but concerns about paternalism have already been raised in the context of health apps<sup>31</sup>.

### Policy challenges and implications

The above distinction between authenticity and agency can be re-captured by two main questions:

1. Does the use of a given AI system lead to the unwarranted distortion of an

individual's beliefs, motivations or decisions?

2. Does the use of a given AI system limit basic liberties or opportunities, or prevent individuals from executing decisions of practical import to their lives?

Answering each question poses additional challenges to developers and policy-makers. Addressing authenticity requires prior deliberation about what conditions need to be fulfilled for external influence to count as (im) permissible. Addressing the external dimension, on the other hand, requires a decision about which options and freedoms are considered essential for autonomy. It also requires deliberation on the permissibility of potential trade-offs between such freedoms.

There exists an extensive body of philosophical literature that is concerned with the first challenge and explicitly lays out what conditions need to be fulfilled for a decision or desire to count as authentic. A prominent approach, developed by Christman, considers a person's decision or desire as authentic if and only if they would not feel alienated from the decision or desire, were they to critically reflect on them<sup>32</sup>. This account emphasises the importance of the individual's point of view when determining whether an external influence counts as undermining autonomy. Coming back to the context of AI, this points towards including users of AI systems much more in the discourse on human autonomy: to determine whether a given system (or the way it is used) is, say, manipulative, it does not suffice to merely observe user behaviour. Instead, we need to test whether users endorse their decisions when given the opportunity to critically reflect on them.

Addressing the second challenge will require explicitly laying out any freedoms, opportunities or decisions that could be affected (positively or negatively; directly or indirectly) by the deployment of any given AI system. Trade-offs should be made explicit, and citizens should be informed about any such limitations or trade-offs.

Identifying potential risks from AI development is a mammoth task. The uncertainty and complexity that surrounds the ethical and social impacts of emerging technologies pose significant challenges to those involved in the governance process. Tackling these challenges requires us to be clear on what it is that we are concerned about

in the first place. Only then can we begin putting into place adequate governance mechanisms to prevent and mitigate potential negative impacts. □

Carina Prunkl  

Institute for Ethics in AI, University of Oxford, Oxford, UK.

 e-mail: carina.prunkl@philosophy.ox.ac.uk

Published online: 23 February 2022  
<https://doi.org/10.1038/s42256-022-00449-9>

### References

1. Raz, J. *The Morality of Freedom* (Clarendon Press, 1986).
2. Korsgaard, C. M., Cohen, G. A., Geuss, R., Nagel, T. Williams, T. & O'Neill, O. *The Sources of Normativity* (Cambridge Univ. Press, 1996).
3. Christman, J. in *The Stanford Encyclopedia of Philosophy* (ed. Zalta, E. N.) (Metaphysics Research Lab, Stanford Univ., 2020); <https://plato.stanford.edu/entries/autonomy-moral/>
4. Roessler, B. *Autonomy: An Essay on the Life Well-Lived* (John Wiley, 2021).
5. Susser, D., Roessler, B. & Nissenbaum, H. *Technology, Autonomy, and Manipulation (Technical Report)* (Social Science Research Network, Rochester, NY, 2019).
6. Kramer, A. D. I., Guillory, J. E. & Hancock, J. T. *Proc. Natl. Acad. Sci. USA* **111**, 8788–8790 (2014).
7. European Commission High-Level Experts Group (HLEG). *Ethics Guidelines for Trustworthy AI (Technical Report B-1049)* (EC, Brussels, 2019).
8. Association for Computing Machinery (ACM). *ACM Code of Ethics and Professional Conduct* (ACM, 2018).
9. Université de Montréal. *Montreal Declaration for a Responsible Development of AI (Forum on the Socially Responsible Development of AI)* (Université de Montréal, 2017).
10. European Committee of the Regions. *White Paper on Artificial Intelligence - A European approach to excellence and trust* (EC, 2020).
11. Organisation for Economic Co-operation and Development. *Recommendation of the Council on Artificial Intelligence (Technical Report OECD/LEGAL/0449)* (OECD 2019); <https://oecd.ai/en/ai-principles>
12. European Commission, Directorate-General for Research and Innovation, European Group on Ethics in Science and New Technologies. *Statement on artificial intelligence, robotics and 'autonomous' systems* (EC, 2018).
13. Floridi, L. & Cowls, J. *Harvard Data Sci. Rev.* **1**, 1–13 (2019).
14. Fjeld, J., Achten, N., Hilligoss, H., Nagy, A. & Srikumar, M. *Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches to Principles for AI (SSRN Scholarly Paper ID 3518482)* (Social Science Research Network, Rochester, NY, 2020); <https://papers.ssrn.com/abstract=3518482>
15. Milano, S., Taddeo, M. & Floridi, L. *Recommender Systems and their Ethical Challenges (SSRN Scholarly Paper ID 3378581)* (Social Science Research Network, Rochester, NY, 2019).
16. Calvo, R. A., Peters, D. & D'Mello, S. *Commun. ACM* **58**, 41–42 (2015).
17. Mik, E. *Law Innov. Technol.* **8**, 1–38 (2016).
18. Helberger, N. *Profiling and Targeting Consumers in the Internet of Things - A New Challenge for Consumer Law (Technical Report)* (Social Science Research Network, Rochester, NY, 2016).
19. Burr, C., Morley, J., Taddeo, M. & Floridi, L. *IEEE Trans. Technol. Soc.* **1**, 21–33 (2020).
20. Morley, J. & Floridi, L. *Sci. Eng. Ethics* **26**, 1159–1183 (2020).
21. Brownsword, R. in *Law, Human Agency and Autonomic Computing* (eds Hildebrandt, M. & Rouvroy, A.) 80–100 (Routledge, 2011).
22. Calvo, R., Peters, D., Vold, K. V. & Ryan, R. in *Ethics of Digital Well-Being* (Philosophical Studies Series, vol. 140) (eds Burr, C. & Floridi, L.) 31–54 (Springer, 2020).
23. Rubel, A., Castro, C. & Pham, A. *Algorithms and Autonomy: The Ethics of Automated Decision Systems* (Cambridge Univ. Press, 2021).
24. Dworkin, G. *The Theory and Practice of Autonomy* (Cambridge Univ. Press, 1988).
25. Mackenzie, C. *Three Dimensions of Autonomy: A Relational Analysis* (Oxford Univ. Press, 2014).
26. Noggle, R. *Am. Philos. Q.* **33**, 43–55 (1996).
27. Elster, J. *Sour Grapes: Studies in the Subversion of Rationality* (Cambridge Univ. Press, 1985).

28. Adomavicius, G., Bockstedt, J. C., Curley, S. P. & Zhang, J. *Info. Syst. Res.* **24**, 956–975 (2013).
29. Ledford, H. *Nature* **574**, 608–609 (2019).
30. Dworkin, G. in *The Stanford Encyclopedia of Philosophy* (ed. Zalta, E. N.) (Metaphysics Research Lab, Stanford Univ. Press, 2020; <https://plato.stanford.edu/archives/fall2020/entries/paternalism/>)
31. Kühler, M. *Bioethics* **36**, 194–200 (2021).

32. Christman, J. *The Politics of Persons: Individual Autonomy and Socio-Historical Selves* (Cambridge Univ. Press, 2009.)

#### Acknowledgements

The author thanks J. Tasioulas, M. Philipps-Brown, C. Veliz, T. Lechterman, A. Dafoe and B. Garfinkel for their helpful comments. Funding: No external funding sources.

#### Competing interests

The author declares no competing interests.

#### Additional information

**Peer review information** *Nature Machine Intelligence* thanks the anonymous reviewers for their contribution to the peer review of this work.