

The Artificial Intelligence of the Ethics of Artificial Intelligence: An Introductory Overview for Law and Regulation

Joanna J. Bryson

July 28, 2019

Abstract

Artificial intelligence (AI) is a technical term often referring to artifacts used to detect contexts for human actions, or sometimes also for machines able to effect actions in response to detected contexts. Our capacity to build such artifacts has been increasing, and with it the impact they have on our society. This does not alter the fundamental roots or motivations of law, regulation, or diplomacy, which rest on persuading humans to behave in a way that provides sustainable security for humans. It does however alter nearly every other aspect of human social behaviour, including making accountability and responsibility potentially easier to trace. This chapter reviews the nature and implications of AI with particular attention to how they impinge on possible applications to and of law.

Keywords: artificial intelligence; regulation; law; diplomacy; accountability; security; social behaviour

For many decades, Artificial Intelligence (AI) has been a schizophrenic field pursuing two different goals: an improved understanding of computer science through the use of the psychological sciences, and an improved understanding of the psychological sciences through the use of computer science. Although apparently orthogonal, these goals have been seen as complementary since progress on one often informs or even progresses the other. Indeed, we have found two factors which have proven to unify the two pursuits. First, the costs of computation and indeed what is actually computable are facts of nature that constrain both natural and artificial intelligence. Second, given the constraints of computability and the costs of computation, greater intelligence relies on the reuse of prior computation. Therefore to

the extent that both natural and artificial intelligence are able to reuse the findings of prior computation, both can be advanced at once.

Neither of the dual pursuits of AI entirely readied researchers for the now glaringly evident ethical importance of the field. Intelligence is a key component of nearly every human social endeavour, and our social endeavours are most of our activities for which we have explicit, conscious awareness. Social endeavours are also the purview of law, and more generally of politics and diplomacy. In short, everything humans deliberately do has been altered by the digital revolution, as well as much of what we do unthinkingly. Often this alteration is in terms of how we can do what we do, for example how we check the spelling of a document, book travel, recall when we last contacted a particular employee, client or, politician, plan our budgets, influence voters from other countries, decide what movie to watch, earn money from performing artistically, discover sexual or life partners, and so on. But what makes the impact ubiquitous is that everything we have done, or chosen not to do, is at least in theory knowable. This fundamentally alters our society because it alters not only what we can do, but how and how well we can know and regulate ourselves and each other.

A great deal has been written about AI ethics recently. Much of it unfortunately has not focussed either on the science of what is computable, nor on the social science of how ready access to more information and more (but mechanical) computational power has altered human lives and behaviour. Rather, a great deal of it has focussed on AI as a thought experiment or ‘intuition pump’ through which we can better understand the human condition or the nature of ethical obligation. This volume is focussed on the law—the day-to-day means by which we regulate our societies and defend our liberties. This chapter sets context for the volume by introducing AI as an applied discipline of science and engineering.

Intelligence is an ordinary process

For the purpose of this introduction I will use an exceedingly-well established definition of intelligence, dating to the seminal monograph on animal behaviour. *Intelligence* is the capacity to do the right thing at the right time. It is the ability to respond to the opportunities and challenges presented by a context. This simple definition is important because it demystifies intelligence, and through it AI. It clarifies both intelligence’s limits and our own social responsibilities, in two ways.

⁰George John Romanes. *Animal intelligence*. London: D. Appleton, 1882.

First, note that intelligence is something that operates at a place and in a moment. It is a special case of *computation*, which is the physical transformation of information¹. Information is not an abstraction². It is physically manifested in light, sound, or materials. Computation and intelligence are therefore also not abstractions. They require time, space, and energy. This is why—when you get down to it—no one is really ever that smart. It is physically impossible to think of everything. We can make tradeoffs: we can for example double the number of computers we use and cut the time of a computation nearly in half. Not quite in half, because there is always an extra cost³ of splitting the task and recombining the outcomes of the processing. But this requires double the space for our two computers, and double the energy in the moment of computation, though the sum of the total energy used is again nearly the same, with the addition of the energy for the overheads. There is no evidence that quantum computing will change this cost equation fundamentally: it should save on not only time but also space, but the energy costs are poorly understood and to date look fiendishly high.

Second, note that the difference between *intelligence* and *artificial intelligence* is only a qualifier. *Artificial* means that something has been made through a human process. This means by default that humans are responsible for it. The artefact actually even more interesting than AI here is a concept: *responsible*. Other animals can be trained to intentionally limit where they place (for example) even the fairly unintentional byproducts of their digestive process, but as far as we know only humans have, can communicate about, and—crucially—can negotiate an explicit concept of responsibility.

Over time, as we recognise more consequences of our actions, our societies tend to give us both responsibility and accountability for these consequences—credit and blame depending on whether the consequences are positive or negative. AI only changes our responsibility as a special case of changing every other part of our social behaviour. Digital technology provides us with *better* capacity to perceive and maintain accounts of actions and consequences, so it should be easier not harder to maintain responsibility and enforce the law. However, whether accountability is easier with AI depends whether and in

¹Michael Sipser. *Introduction to the Theory of Computation*. Second. Boston, MA: PWS, Thompson, 2005.

²Claude Elwood Shannon. “A mathematical theory of communication”. In: *Bell system technical journal* 27.3 (1948), pp. 379–423.

³an overhead, cf. Ajay D Kshemkalyani and Mukesh Singhal. *Distributed computing: principles, algorithms, and systems*. Cambridge University Press, 2011.

what ways we deploy the capacities digital technology affords. Without care and proper measures, the increased capacity for communication that information communication technology (ICT) provides may be used to diffuse or obscure responsibility. One solution is to recognise in law that the lack of such care and measures for promoting accountability in processes concerning digital artefacts is a form of negligence. Similarly, we could declare unnecessary obfuscation of public or commercial processes a deliberate and culpable evasion of responsibility.

Note that the simplicity of the definitions introduced in this section is extremely important as we move towards law and regulation of systems and societies infused with AI. In order to evade regulation or responsibility, the definition of intelligence is often complicated in manifestos by notions such as sentience, consciousness, intentionality and so forth. I will return to these issues below, but what is essential when considering AI in the context of law is the understanding that no fact of either biology (the study of life) nor computer science (the study of what is computable) names a necessary point at which human responsibility should end. Responsibility is not a fact of nature. Rather, the problem of governance is as always to design our artefacts—including the law itself—in a way that helps us maintain enough social order so that we can sustain human flourishing.

AI—including machine learning—occurs by design

AI only occurs by and with design. AI is only produced intentionally, for a purpose, by one or more members of human society. That act of production requires design decisions concerning at a minimum the information input to and output from the system, and also where and how the computation required to transform that information will be run. These decisions entail also considerations of energy consumption and time that can be taken in producing as good a system as possible. Finally, any such system can and should be defended with levels of both cyber and physical security appropriate to the value of the data transmitted or retained, and the physical capacities of the system if it acts on the world⁴.

The tautology that AI is always generated by design extends to *machine learning* (ML), which is one means of developing AI wherein computation is used to discover useful regularities in data. Systems can then be built

⁴Note that these observations show how under-informed about basic systems engineering the idea of a machine converting the world into paperclips is, as per Nick Bostrom. *Superintelligence: Paths, dangers, strategies*. Oxford University Press, 2014, p. .

to exploit these regularities either to categorize, make predictions, or select actions directly. The mere fact that part of the process of design has been automated does not mean that the system itself is not designed. The choice of ML algorithm, the data fed into it to train it, the point at which it is considered adequately trained to be released, how that point is detected by testing, whether that testing is ongoing if the learning continues during the system’s operation—all of these things are design decisions that not only must be made, but that also can easily be documented. As such, any individual or organisation that produces AI could always be held to account by being asked to produce documentation of these processes.

Despite the fact documentation of such decisions and records of testing outcomes are easy to produce, good practice is not always followed⁵. This is as much a matter for the law as any other sloppy or inadequate manufacturing technique⁶. What process is deemed adequate for commercial products or even private enjoyment is determined by some combination of expertise and precedent. Whether these processes have been followed *and* documented can easily be checked either before a product is licensed, after a complaint has been made, or as a part of routine inspection.

Although actual algorithms are abstractions, that only means algorithms in themselves are not AI. In computer science, an algorithm is just a list of instructions to be followed, like a recipe in baking⁷. Just as a strand of DNA in itself is not life—it has no capacity to reproduce itself—so instruction sets require not only input (data) but also physical computation. Without significant complex physical infrastructure to execute their instructions, both DNA and AI algorithms are vacuous. The globally-largest technology corporations have almost inconceivably vast infrastructure for every aspect of storing, processing, and transmitting the information that is their business. This infrastructure includes means to generate electric power and provide secure communication as well as means to do computation.

These few leading corporations further provide these capacities also as service infrastructure to a significant percentage of the world’s other ICT companies—of course, at a cost. The European Union (EU) has committed to investing substantial public resources in developing a localised equiva-

⁵Michael Hüttermann. *DevOps for Developers*. Apress, Springer, 2012.

⁶Joshua A. Kroll et al. “Accountable Algorithms”. In: *University of Pennsylvania Law Review* 165 (2017), pp. 633–706.

⁷The term algorithm is currently often misused to mean an AI system by those unclear on the distinction between design, program, data, and physical computing system e.g. Cathy O’Neil. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Broadway Books, 2016, p. .

lent of this computational infrastructure resource, as they have previously done with both commercial aviation and global positioning systems. The EU may also attempt to build a parallel data resource, though this is more controversial. There has also been some discussion of ‘nationalising’ significant technology infrastructure, though that idea is problematic given that the Internet is transnational. *Transnationalising technology ‘giants’* is also discussed further below.

Digital technology empowers us to do all sorts of things, including obfuscating or simply deleting records or the control systems they refer to. We can make systems either harder or easier to understand using AI⁸. These are design decisions. The extent to which transparency and accountability should be required in legal products is also a design decision, though here it is legislators, courts, and regulators that design a regulatory framework. What is important to realize here is that it is perfectly possible to mandate that technology be designed to comply with laws, including any ensuring traceability and accountability of the human actions involved in the design, running, and maintenance of intelligent systems. In fact, given that the limits of ‘machine nature’ are far more plastic than those of human nature, it is far more sensible to minimise the amount of change to laws and rather maximise the extent of required compliance to and facilitation of extant laws⁹.

The performance of designed artefacts is readily explainable

Perhaps in the desire to evade either laws of nations or the above-mentioned laws of nature, many deeply respected AI professionals have claimed that the most promising aspects of AI would be compromised if AI were to be regulated¹⁰. For example, the idea that maintaining standard rights to

⁸Kroll et al., “Accountable Algorithms”.

⁹Joanna J. Bryson, Mihailis E. Diamantis, and Thomas D. Grant. “Of, for, and by the people: the legal lacuna of synthetic persons”. In: *Artificial Intelligence and Law* 25.3 (Sept. 2017), pp. 273–291. ISSN: 1572-8382. DOI: 10.1007/s10506-017-9214-9. URL: <https://doi.org/10.1007/s10506-017-9214-9>; Margaret Boden et al. *Principles of Robotics*. The United Kingdom’s Engineering and Physical Sciences Research Council (EPSRC). Apr. 2011. URL: <https://www.epsrc.ac.uk/research/ourportfolio/themes/%20engineering/activities/principlesofrobotics/>.

¹⁰My assertion about the ‘deeply respected’ relates to claims I’ve heard in high-level policy settings, but haven’t been able to find in print. However, for examples of the rhetoric see Cassie Kozyrkov. “Explainable AI won’t deliver. Here’s why.” In: *Hackernoon* (Nov.

explanation—demonstration of due process—would eliminate the utilisation of many advanced machine learning techniques, because these are too complex for their exact workings to be knowable. This last sort of assertion fails to take into account the present standards for accountability in corporate law. If a company is audited, that audit never reaches to explaining the workings of the brain synapses or gene regulation of that company’s employees. Rather, we look for audit trails—or perhaps witnesses—indicating that humans have followed appropriate procedures.

AI may reduce the number of people who can be put on a witness stand to describe their recollections of events or motivations, but it enables a standard of record keeping that would be unbearably tedious in non-digital processes. It is not the case that all AI systems are programmed to keep such records, nor that all such records are maintained indefinitely. But it *is* the case that *any* AI system can be programmed for this, and programmed using good systems of logging of the design, development, training, testing, as well as the operation of the systems. Further, individuals or institutions can choose how, where, and how long to store this logging data. Again, these are design decisions for both AI systems and the institutions that create them. There are available standards for adequate logging for generating proof of due diligence or even explanation of behaviour is. Norms of use for these standards can be set and enforced¹¹.

What matters for human justice is that humans do the right things. We do not need to check exactly how a machine learning algorithm works any more than we need to completely understand the physics of torque to regulate bicycle riding in traffic. Our concerns with AI are that it is used in a way that is lawful. We want to know for example that products comply to their claims, that individual users are not spied upon or unfairly disadvantaged, that foreign agencies were not able to illicitly insert false information into training a machine learning data set or into a newsfeed.

All AI affords the possibility of maintaining precise accounts of when, how, by whom, and with what motivation it has been constructed. Indeed, this is true of artefacts in general, but digital artefacts are particularly amenable to automating the process. The very tools used to build the sys-

2018). URL: <https://hackernoon.com/%20explainable-ai-wont-deliver-here-s-why-6738f54216be>; Erdem. “The trade-off in machine learning: Accuracy vs explainability”. In: *Medium* (Dec. 2018). URL: <https://medium.com/@erdemkalayci/%20the-tradeoff-in-machine-learning-accuracy-vs-explainability-fbb13914fde2>.

¹¹Joanna J. Bryson and Alan F. T. Winfield. “Standardizing Ethical Design for Artificial Intelligence and Autonomous Systems”. In: *Computer* 50.5 (May 2017), pp. 116–119. ISSN: 0018-9162. DOI: 10.1109/MC.2017.154.

tem can also be set to capture and prompt for this kind of information. We can similarly track the construction, application, and outcomes of any validating tests. Further, even the most obscure AI system after development can be treated entirely as a blackbox and still tested to see what variation in inputs varies the outputs¹². Even where performance is stochastic, statistics can tell us the probability of various outcomes, again a sort of information to which the law is already accustomed, for example in medical outcomes.

In fact though almost no AI systems are entirely opaque. Systems with AI are generally far less opaque than human reasoning and less complex than a government or ecosystem. There is a decades-old science of examining complex models by using simpler ones, which is already accelerating to serve the sectors that are already well regulated which are of course (like all sectors) increasingly using AI¹³. And of course many forms of AI, built either with or without the use of ML, readily produce explanations themselves¹⁴.

To return to one of the assertions at the beginning of this section, it is also wrong to assume that AI is not already regulated. All human activity, particularly commercial activity, occurs in the context of some sort of regulatory framework¹⁵. The question is how to continue to optimise this framework in light of the changes to society and its capacities introduced by AI and ICT more generally.

¹²This process is coming to be called (as of this writing) forensic analysis, see e.g. Joseph R. Barr and Joseph Cavanaugh. “Forensics: Assessing model goodness: A machine learning view”. In: *Robotic Intelligence*. 2019, pp. 17–23. DOI: 10.1142/9789811203480_0003. URL: https://www.worldscientific.com/doi/abs/10.1142/%209789811203480_0003.

¹³Patrick Hall. “On the Art and Science of Machine Learning Explanations”. In: *arXiv preprint arXiv:1810.02909* (2018).

¹⁴Stephen Cranefield et al. “No Pizza for You: Value-based Plan Selection in BDI Agents”. In: *IJCAI Proceedings*. Ed. by Carles Sierra. Melbourne, Aug. 2017; Jiaming Zeng, Berk Ustun, and Cynthia Rudin. “Interpretable classification models for recidivism prediction”. In: *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 180.3 (2017), pp. 689–722. ISSN: 1467-985X. DOI: 10.1111/rssa.12227. URL: <http://dx.doi.org/10.1111/rssa.12227>.

¹⁵Miles Brundage and Joanna J. Bryson. *Smart Policies for Artificial Intelligence*. in preparation, available as arXiv:1608.08196. 2017.

Intelligence increases by exploiting prior computation

The fact that computation is a physical process limits how much can be done *de novo* in the instant during which intelligence must be expressed—when action must be taken to save a system from a threat or to empower it through an opportunity. For this reason, much of intelligence exploits computation already done, or rather artefacts produced that preserve the outcomes of that computation. We can understand the outcomes not only of culture but of biology this way. It is not only that organisms can only exploit opportunities they can perceive, it is also that they tend only to perceive what they are equipped to exploit—both capacities for perception and action evolve together. Similarly, culture passes us the tools that others have not only invented but, of all those inventions, the ones that produce the greatest impact relative the costs of transmission, where such costs include both time (suggesting missed opportunities) and the likelihood of adequately faithful replication¹⁶. Culture itself evolves because that gives us and it efficacy¹⁷.

Much of the recent immense growth of AI has been due to improved capacities to ‘mine’ using ML the existing discoveries of humanity and nature more generally¹⁸. The result of this is of course that the good comes with the bad. We mine not only knowledge but stereotypes—and if we allow AI to take action, prejudice—when we mine human culture¹⁹. This is not

¹⁶Ivana Čaće and Joanna J. Bryson. “Agent Based Modelling of Communication Costs: Why Information can be Free”. In: *Emergence and Evolution of Linguistic Communication*. Ed. by C. Lyon, C. L. Nehaniv, and A. Cangelosi. London: Springer, 2007, pp. 305–322; Kenny Smith and Elizabeth Wonnacott. “Eliminating unpredictable variation through iterated learning”. In: *Cognition* 116.3 (2010), pp. 444–449. ISSN: 0010-0277. DOI: <https://doi.org/10.1016/j.cognition.2010.06.004>. URL: <http://www.sciencedirect.com/science/article/pii/S0010027710001320>.

¹⁷Alex Mesoudi, Andrew Whiten, and Kevin N. Laland. “Towards a unified science of cultural evolution”. In: *Behavioral and Brain Sciences* 29.4 (2006), pp. 329–347. DOI: 10.1017/S0140525X06009083; Joanna J. Bryson. “Embodiment versus Memetics”. In: *Mind & Society* 7.1 (June 2008), pp. 77–94; Joanna J. Bryson. “Artificial Intelligence and Pro-Social Behaviour”. In: *Collective Agency and Cooperation in Natural and Artificial Systems: Explanation, Implementation and Simulation*. Ed. by Catrin Misselhorn. Vol. 122. Philosophical Studies. Berlin: Springer, Oct. 2015, pp. 281–306; Daniel C. Dennett. *From Bacteria to Bach and Back*. Allen Lane, 2017.

¹⁸Thomas B Moeslund and Erik Granum. “A survey of computer vision-based human motion capture”. In: *Computer vision and image understanding* 81.3 (2001), pp. 231–268; Sylvain Calinon et al. “Learning and reproduction of gestures by imitation”. In: *IEEE Robotics & Automation Magazine* 17.2 (2010), pp. 44–54.

¹⁹Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. “Semantics derived automatically from language corpora contain human-like biases”. In: *Science* 356.6334

a special feature of AI, as mentioned above this is how nature works as well²⁰. Further, show that at least some of what we call ‘stereotype’ reflects aspects of present-day conditions, such as what proportion of jobholders for a particular position have a particular gender. Thus some things we have agreed are bad (e.g. that it is sexist to expect programmers to be male) are aspects of our present culture we have therefore implicitly agreed we wish to change.

One theory for explaining the explosion in what we recognise as AI (that is, of AI with rich, demonstrably human-like, and previously human-specific capacities like speech production and comprehension, or face recognition) is less a consequence of new algorithms than of new troves of data and increased computation speeds. We can expect such explosions of capacities based on the strategy of mining past solutions to soon plateau, when artificial and human intelligence come to be sharing nearly the same, though still-expanding boundary of extant knowledge. In fact, we can also expect this boundary to be expanding faster now, given the extra computational resources we are bringing not only through digital hardware, but by increasing access to other human minds. For humanity, ICT reduces the aforementioned overheads of combining concurrent search. We all get smarter as our culture expands to embrace more—and more diverse—minds²¹. However, the fact that we can exploit our own computation to build AI, or that we can build our own native as well as systemic intelligence by using AI, does not mean that we are replaceable with AI. As will be explained in the next sections, AI cannot be used to replicate humans, and this has substantial consequences for law and regulation.

(2017), pp. 183–186. ISSN: 0036-8075. DOI: 10.1126/science.aal4230. URL: <http://science.sciencemag.org/content/356/6334/183>.

²⁰Molly Lewis and Gary Lupyan. “Language use shapes cultural norms: Large scale evidence from gender”. In: *Proceedings of the Annual Meeting of the Cognitive Science Society (CogSci)*. also in prep. for journal publication. Madison, WI, 2018, pp. 2041–2046.

²⁰Caliskan, Bryson, and Narayanan, “Semantics derived automatically from language corpora contain human-like biases”.

²¹Anita Williams Woolley et al. “Evidence for a Collective Intelligence Factor in the Performance of Human Groups”. In: *Science* 330.6004 (29 October 2010), pp. 686–688; Barton H. Hamilton, Jack A. Nickerson, and Hideo Owan. “Diversity and Productivity in Production Teams”. In: *Advances in the Economic Analysis of Participatory and Labor-Managed Firms*. 2012, pp. 99–138. DOI: 10.1108/S0885-3339(2012)0000013009. URL: <https://www.emeraldinsight.com/doi/abs/10.1108/%20S0885-3339%282012%290000013009>; Feng Shi et al. “The wisdom of polarized crowds”. In: *Nature Human Behaviour* 3 (2019), pp. 329–336.

AI cannot produce fully replicated humans (all models are wrong)

When computer science is mistaken for a branch of mathematics, many important implications of computation being a physical process are lost. For example, AI is wrongly perceived as a path towards human immortality. First, the potential of ‘uploading’ human intelligence in any meaningful sense is highly dubious. Technologically, brains cannot be ‘scanned’ and replicated in any other material than another brain, as their computational properties depend on trillions of temporal minutiae²². Creating a second, identical human to host that new brain is not only physically intractable, but would be cloning—both unethical and illegal, at least in the European Union. Second, even if²³ we could rather somehow upload adequate abstractions of our own minds, we should not confuse this with actually having spawned a digital replica. A digital clone might be of use for example to offload canned email replies²⁴, or to create somewhat interactive interfaces for historical story telling²⁵.

Many have argued that the moral intuitions, motivations, even the aesthetics of an enculturated ape can in no way be meaningfully embedded in a device that shares nothing of our embodied physical (‘phenomenological’) experience²⁶. Nothing we build from metal and silicon will ever share our phenomenology as much as a rat or cow, and few see cows or rats as viable vessels of our posterity. Yet whether such digital artefacts are viewed as adequate substitutes for a real person depends on what one values about that person. For example, if one is valuing one’s own capacity to control

²²Yoonsuck Choe, Jaeroek Kwon, and Ji Ryang Chung. “Time, Consciousness, and Mind Uploading”. In: *International Journal of Machine Consciousness* 04.01 (2012), pp. 257–274. DOI: 10.1142/S179384301240015X. URL: <http://www.worldscientific.com/doi/abs/10.1142/%20S179384301240015X>.

²³as some would suggest, see Murray Shanahan. *The technological singularity*. MIT Press, 2015, for a review.

²⁴Mark Dredze et al. “Intelligent email: Reply and attachment prediction”. In: *Proceedings of the 13th international conference on Intelligent user interfaces*. ACM, 2008, pp. 321–324.

²⁵David Traum et al. “New Dimensions in Testimony: Digitally preserving a Holocaust survivor’s interactive storytelling”. In: *Proceedings of the Eighth International Conference on Interactive Digital Storytelling*, pp. 269–281.

²⁶Frank Pasquale. “Two concepts of immortality: Reframing public debate on stem-cell research”. In: *Yale Journal of Law and the Humanities* 14 (2002), pp. 73–121; Bryson, “Embodiment versus Memetics”; Guy Claxton. *Intelligence in the flesh: Why your mind needs your body much more than it thinks*. Yale University Press, 2015; Dennett, *From Bacteria to Bach and Back*.

the lives of others, many already turn to the simple technology of a will to control sometimes quite intimate aspects of the lives of those chosen to be their heirs. Thus it seems clear that there will be those who spend millions or even billions of dollars, euro, or rubles on producing digital clones they are literally deeply invested in believing to be themselves²⁷.

Even if we could somehow replicate ourselves in an artefact, the mean time to obsolescence of digital technologies and formats is far, far shorter than the average human life expectancy, presently nearing ninety years. This quick obsolescence is true not only of our physical technology, but also of our fashion. Unquestionably any abstracted digital self-portrait would follow fashion in reflecting an aspect of our complex selves that will have been culturally-appropriate only in a specific moment. It would not be possible from such an abstraction to fully model how our own rich individual being would have progressed through time, let alone through biological generations. Such complete modelling opposes the meaning of *abstraction*. An unabstracted model would again require biological cloning, and even there after many generations would fall out of ecological ‘fashion’ or appropriateness as evolution takes its course.

With apologies to both Eisenhower and, all abstractions are wrong, but producing abstractions is essential. By the definition used in this chapter, intelligent action is an abstraction of the present context, therefore producing an abstraction is the essence of intelligence. But that abstraction is only a snapshot of the organism, it is not the organism itself.

Reproducing our full organism is not required for many aspects of what calls ‘positive immortality’. Replicating our full selves is certainly not essential to writing fiction or otherwise making a lasting contribution to a culture or society, or an irrevocable impact on the ecosystem. But the purpose of this chapter is to introduce AI from the perspective of maintaining social order — that is, from the perspective of law and regulation. As will be discussed below, the methods for enforcing law and regulation are founded on the evolved priorities of social animals. Therefore any intelligent artefacts representing such highly abstracted versions of an individual human

²⁷Pasquale, “Two concepts of immortality: Reframing public debate on stem-cell research”, questions such expenditures, or even those of in vitro fertilisation, on the grounds economic fairness.

²⁷G. E. P. Box. “Robustness in the strategy of scientific model building”. In: *Robustness in statistics*. Ed. by R. L. Launer and G. N. Wilkinson. New York, NY: Academic Press, 1979, pp. 201–236.

²⁷Pasquale, “Two concepts of immortality: Reframing public debate on stem-cell research”.

are not relevant to the law except perhaps as the intellectual property of their creator.

AI itself cannot be dissuaded by law or treaty

There is no way to ensure that an artefact²⁸ could be held legally accountable. Many people think the purpose of the law is to compensate, and obviously if we allow a machine to own property or at least wealth then it could in some sense compensate for its errors or misfortune. However, the law is really primarily designed to maintain social order by dissuading people from doing wrong. Law dissuades by making it clear first what actions are considered wrong, and second what are the costs and penalties for committing these wrong acts. This is even more true of policy and treaties, which are often constructed after long periods of negotiated agreement between peers or at least sufficiently powerful actors about what these wrongs and costs are. The Iran Nuclear Deal²⁹ is an excellent example of this.

Of course all of these systems of governance can also generate revenue, which may be used by governments to some extent to right wrongs. However, none of the costs or penalties that courts can impose will matter to an AI system. While we can easily write a program that says “Don’t put me in jail!” we cannot program the full, systemic aversion to the loss of social status and years of one’s short life that the vast majority of humans experience by birthright. In fact, not only humans but many social species find isolation and confinement deeply aversive—guppies can die of fright if separated from their school, and factory farming has been shown to drive pigs to exhibit symptoms of severe mental illness³⁰.

We might add a bomb, camera, and timer to a robot, and program the bomb to destruct if the camera has seen no humans (or other robots) for ten minutes. Reasoning by empathy you might think this machine is far more dissuaded than a human, who can easily spend more than ten minutes alone. But empathy is a terrible system for establishing universal ethics—it

²⁸With no human components Christian List and Philip Pettit. *Group agency: The possibility, design, and status of corporate agents*. Oxford University Press, 2011.

²⁹Kenneth Katzman and Paul K Kerr. *Iran nuclear agreement*. Tech. rep. R43333. Library of Congress, Congressional Research Service, May 2016. URL: www.crs.gov.

³⁰Françoise Wemelsfelder. “The scientific validity of subjective concepts in models of animal welfare”. In: *Applied Animal Behaviour Science* 53.1 (1997). Basic and Applied Aspects of Motivation and Cognition, pp. 75–88. ISSN: 0168-1591. DOI: [https://doi.org/10.1016/S0168-1591\(96\)01152-5](https://doi.org/10.1016/S0168-1591(96)01152-5). URL: <http://www.sciencedirect.com/science/article/pii/S0168159196011525>.

works best on those most like you³¹. The robot’s behaviour could easily be utterly unaltered by this contrivance, and so it could not be said to suffer at all by the technical definitions of suffering³². Even if the robot could detect and reason about the consequences of its new situation, it would not feel fear, panic, or any other systemic aversion, although depending on its goals it may alter its planning to favor shorter planning horizons.

Law has been invented by—we might even say ‘coevolved with’—our societies to hold humans accountable, thus only humans can be held accountable with it. Even the extension of legal personality to corporations only works to the extent that real humans who have real control over those corporations suffer if the corporation does wrong. The overextension of legal personhood to corporations that are designed to fail is called making a ‘shell company’. Similarly, if you build an AI system and allow it to operate autonomously, it is essential that the person who chooses to allow the system to operate autonomously is the one who will go to jail, be fined, etc. if the AI system transgresses the law. There is no way to make the AI system itself accountable. AI being itself held accountable would be the ultimate shell company³³.

The implicit principles that underly our capacity to coordinate and cooperate through the law and its dissuasions have also coevolved with our advanced societies. We share many of our cognitive attributes—including perception and action capacities, and importantly, motivations—with other apes. Yet we also have specialist motivations and capacities reflecting our highly social nature³⁴. No amount of intelligence in itself necessitates social competitiveness, neither does it demand the desire to be accepted by an ingroup, to dominate an outgroup, nor to achieve recognition within an ingroup. These are motivations that underlie human cooperation and

³¹Paul Bloom. *Against empathy: The case for rational compassion*. Random House, 2017.

³²Wemelsfelder, “The scientific validity of subjective concepts in models of animal welfare”; Daniel C. Dennett. “Why You Can’t Make a Computer that Feels Pain”. In: *Brainstorms*. page numbers are from the 1986 Harvester Press Edition, Brighton, Sussex. Montgomery, Vermont: Bradford Books, 1978, pp. 190–229; Bryson, “Artificial Intelligence and Pro-Social Behaviour”; Margaret A. Boden. “Robot says: Whatever (The robots won’t take over because they couldn’t care less)”. In: *Aeon* (13 August 2018). originally a lecture at the Leerhulme Centre for the Future of Intelligence. URL: <https://aeon.co/essays/%20the-robots-wont-take-over-because-they-couldnt-care-less>.

³³Bryson, Diamantis, and Grant, “Of, for, and by the people: the legal lacuna of synthetic persons”.

³⁴David Michael Stoddart. *The Scented Ape: The Biology and Culture of Human Odour*. Cambridge University Press, Nov. 1990.

competition that result from our evolutionary history³⁵. None of this is necessary—and much of it is even incoherent—from the perspective of an artefact. Artefacts are definitionally designed by human intent, not directly by evolution. With these intentional acts of authored human creation³⁶ comes not only human responsibility, but an entirely different landscape of potential rewards and design constraints³⁷.

AI and ICT impact every human endeavour

Given that AI can always be built to be explainable, and that only humans can be held to account, assertions that AI itself should be trustworthy, accountable, or responsible are completely misguided. If only humans can be held to account; from a legal perspective the goal for AI transparency is to ensure that human blame can be correctly apportioned. Of course there are also other sorts of transparency such as those that support ordinary users in establishing the correct boundaries they have with their systems (defending their own interests), and for developers or other practitioners to be able to debug or customize an AI system³⁸. AI can be reliable but not trustworthy—it should not require a social compact or leap of faith³⁹. Consumers and governments alike should have confidence that there is a fact

³⁵Stoddart, *The Scented Ape: The Biology and Culture of Human Odour*; Ruth Mace. “The co-evolution of human fertility and wealth inheritance strategies”. In: *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 353.1367 (1998), pp. 389–397. ISSN: 0962-8436. DOI: 10.1098/rstb.1998.0217. URL: <http://rstb.royalsocietypublishing.org/content/353/1367/%20389>; Jillian J. Jordan et al. “Uncalculating cooperation is used to signal trustworthiness”. In: *Proceedings of the National Academy of Sciences* (2016). DOI: 10.1073/pnas.1601280113. URL: <http://www.pnas.org/content/early/2016/07/19/%201601280113.abstract>; Simon T. Powers, Carel P. van Schaik, and Laurent Lehmann. “How institutions shaped the last major evolutionary transition to large-scale human societies”. In: *Philosophical Transactions of the Royal Society B: Biological Sciences* 371.1687 (2016), p. 20150098. DOI: 10.1098/rstb.2015.0098. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/%20rstb.2015.0098>.

³⁶The choice to create life through childbirth is not the same. While we may author some of childrearing, the dispositions just discussed are shared with other primates, and are not options left to parents or other conspecifics to determine.

³⁷cf. Joanna J. Bryson. “Patience is not a virtue: the design of intelligent systems and systems of ethics”. In: *Ethics and Information Technology* 20.1 (Mar. 2018), pp. 15–26. ISSN: 1572-8439. DOI: 10.1007/s10676-018-9448-6. URL: <https://doi.org/10.1007/s10676-018-9448-6>.

³⁸Bryson and Winfield, “Standardizing Ethical Design for Artificial Intelligence and Autonomous Systems”.

³⁹Onora O’Neill. *A question of trust: The BBC Reith Lectures 2002*. Cambridge University Press, 2002.

of the matter that they can determine at will about who is responsible for the AI infused systems we incorporate into our homes, our business processes, and our security.

Every task we apply our conscious minds to—and a great deal we do implicitly—we do using our intelligence. Artificial Intelligence therefore can affect everything we are aware of doing and a great deal we have always done without intent. As mentioned earlier, using even fairly trivial and ubiquitous AI we recently demonstrated that human language contains our implicit biases, and further those in many cases reflect our lived realities⁴⁰. In reusing and reframing our previous computation, AI allows us to see truths we hadn't previously known about ourselves, including how we transmit stereotypes⁴¹, but it doesn't automatically or magically improve on us without effort. Caliskan, Bryson, and Narayanan show also that the outcome of the famous study showing that given otherwise-identical resumé's, individuals with stereotypically African American names were half as likely to be invited to a job interview as individuals with European American names. Smart corporations are now using carefully-programmed AI to avoid the implicit biases at the early stages of their HR processes to get the sort of diverse CVs to the short-list stage where hiring decisions can—with explicit care and intention—avoid perpetuating the mistakes of our past.

The idea of having 'autonomous' AI systems 'value aligned' is therefore also likely also misguided. While it is certainly necessary to acknowledge and understand the extent to which implicit values and expectations must be embedded in any artefact⁴², designing for such embedding is not sufficient to create a system that is autonomously moral. Indeed, if as I have argued a system cannot be accountable, it may also not in itself be held as a moral agent. The issue should not be 'embedding' our intended (or asserted) values in our machines, but rather ensuring that our machines can

⁴⁰Caliskan, Bryson, and Narayanan, "Semantics derived automatically from language corpora contain human-like biases".

⁴¹Lewis and Lupyán, "Language use shapes cultural norms: Large scale evidence from gender".

⁴²Marianne Bertrand and Sendhil Mullainathan. "Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination". In: *The American Economic Review* 94.4 (2004), pp. 991–1013.

⁴²Jeroen van den Hoven. "ICT and Value Sensitive Design". In: *The Information Society: Innovation, Legitimacy, Ethics and Democracy In honor of Professor Jacques Berleur s.j.* Ed. by Philippe Goujon et al. Boston, MA: Springer US, 2007, pp. 67–72. ISBN: 978-0-387-72381-5; Aimee van Wynsberghe. "Designing Robots for Care: Care Centered Value-Sensitive Design". In: *Science and Engineering Ethics* 19.2 (June 2013), pp. 407–433. ISSN: 1471-5546. DOI: 10.1007/s11948-011-9343-6. URL: <https://doi.org/10.1007/s11948-011-9343-6>.

allow the expression of the mutable intentions of ourselves, their operators.

Only through correctly expressing our intentions should AI incidentally telegraph our values. Individual liberty, including freedom of opinion and of thought, are absolutely critical not only to human well being but to a robust and creative society⁴³. Allowing ‘values’ to be enforced by the enfold-ing curtains of interconnected technology invites gross excesses by powerful actors against those they consider vulnerable, a threat, or just unimportant⁴⁴. Even supposing a power that is demonstrably benign, allowing it the mechanisms for technological autocracy creates a niche that may facilitate a less-benign power— whether through change of hands, corruption of the original power, or just corruption of the systems communicating its will. Finally, who or what is a powerful actor is also altered by ICT, where clandestine networks can assemble—or be assembled⁴⁵—out of small numbers of anonymous individuals acting in a well coordinated way, even across borders.

Theoretical biology tells us that where there is greater communication, there is a higher probability of cooperation⁴⁶. *Cooperation* has nearly entirely positive connotations, but it is in many senses almost neutral—nearly all human endeavours involve cooperation, and while these generally benefit many humans, some are destructive to many others. Further, the essence of cooperation is moving some portion of autonomy from the individual to a group⁴⁷. The extent of autonomy an entity has is the extent to which it determines its own actions⁴⁸. Individual and group autonomy must to some

⁴³Julie E. Cohen. “What privacy is for”. In: *Harvard Law Review* 126 (May 2013), pp. 1904–1933.

⁴⁴Brett Frischmann and Evan Selinger. *Re-engineering humanity*. Cambridge University Press, 2018; Miles Brundage et al. *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation*. Tech. rep. <https://maliciousaireport.com/>. Future of Humanity Institute, University of Oxford, Centre for the Study of Existential Risk, University of Cambridge, Center for a New American Security, Electronic Frontier Foundation, and OpenAI, Feb. 2018.

⁴⁵Carole Cadwalladr. “I made Steve Bannon’s psychological warfare tool’: meet the data war whistleblower”. In: *The Observer* (18 March 2018).

⁴⁶Joan Roughgarden, Meeko Oishi, and Erol Akçay. “Reproductive Social Behavior: Cooperative Games to Replace Sexual Selection”. In: *Science* 311.5763 (2006), pp. 965–969. DOI: [10.1126/science.1110105](https://doi.org/10.1126/science.1110105). URL: <http://www.sciencemag.org/content/311/5763/965.abstract>.

⁴⁷Bryson, “Artificial Intelligence and Pro-Social Behaviour”.

⁴⁸Harvey Armstrong and Robert Read. “Western European micro-states and EU autonomous regions: The advantages of size and sovereignty”. In: *World Development* 23.7 (1995), pp. 1229–1245. ISSN: 0305-750X. DOI: [http://dx.doi.org/10.1016/0305-750X\(95\)00040-J](http://dx.doi.org/10.1016/0305-750X(95)00040-J). URL: <http://www.sciencedirect.com/science/article/pii/S0305750X9500040J>; Maeve Cooke. “A space of one’s own: Autonomy, privacy,

extent trade off, though there are means of organising groups that offer more or less liberty for their constituent parts.

Many people are (falsely) preaching that ML is the new AI, and (again falsely) that the more data ML is trained on, the smarter the AI. ML is actually a statistical process we use for programming some aspects of AI. Thinking that bigger data is necessarily better begs the question, better for what? Basic statistics teaches us that the number of data points we need to make a prediction is limited by the amount of variation in that data, providing only that the data is a true random sample of its population⁴⁹. So there are natural limits for any particular task on how much data is needed—except perhaps for surveillance. What we need for science or medicine may require only a minuscule fraction of a population. However, if we want to spot specific individuals to be controlled, dissuaded, or even promoted, then of course we want to “know all the things.”⁵⁰

The changing costs and benefits of investment at the group level that Roughgarden, Oishi, and Akçay describe has other consequences beyond privacy and liberty. ICT facilitates blurring the distinction between customer and corporation, or even the definition of an economic transaction. Customers now do real labour for the corporations to whom we give our custom: pricing and bagging groceries, punching data at ATMs for banks, filling in forms for airlines and so forth⁵¹. The value of this labour is not directly remunerated—we assume that we receive cheaper products in return, and as such our loss of agency to these corporations might be seen as a form of bartering. ‘Free’ services like search and email may better be understood as information bartering⁵². These transactions are not denominated with a price, which means that ICT facilitates a black or at least opaque market reducing both measured custom and therefore tax revenue. This is true for everyone who uses Internet services and interfaces, even ignoring the present controversies over definitions of employment raised by platforms⁵³. Our in-

liberty”. In: *Philosophy & Social Criticism* 25.1 (1999), pp. 22–53. DOI: 10.1177/019145379902500102. URL: <http://dx.doi.org/10.1177/019145379902500102>.

⁴⁹Meng18.

⁵⁰Mark Andrejevic. “Automating Surveillance”. In: *Surveillance & Society* 17.1/2 (2019), pp. 7–13.

⁵¹Bryson, “Artificial Intelligence and Pro-Social Behaviour”.

⁵²Joanna J. Bryson. “The Past Decade and Future of AI’s Impact on Society”. In: *Towards a New Enlightenment? A Transcendent Decade*. OpenMind BBVA. commissioned, based on a white paper also commissioned, that by the OECD. Madrid: Taylor, Mar. 2019. URL: <https://www.bbvaopenmind.com/en/articles/%20the-past-decade-and-future-of-ais-impact-on-society/>.

⁵³though see Tim O’Reilly. *WTF? What’s the Future and why It’s Up to Us*. New York:

ability to assign value to these transactions may also explain the mystery of why AI doesn't seem to be increasing productivity⁵⁴.

AI then gives us new ways to do everything we do intentionally and a great deal else. The extent to which AI makes different tasks easier and harder varies in ways that are not intuitive. This also increases and decreases the values of human skills, knowledge, social networks, personality traits, and even locations, and alter the calculations of identity and security. Fortunately, AI also gives us tools for reasoning and communicating about all these changes, and for adjusting to them. But this makes group-level identity itself more fluid, complicating our ability to govern.

Who's in charge? AI and governance

Despite all of this fluctuation, there are certain things that are invariant to the extent of computational resource and communicative capacities. The basic nature of humans as animals of a certain size and metabolic cost, and the basic drives that determine what gives us pleasure, pain, stress, and engagement are not altered much. How we live is and always will be enormously impacted by how our neighbours live, as we share geographically-related decisions concerning investment in air, water, education, health, and security. For this reason there will always be some kind of geography-based governance. The fundamental ethical framework we have been negotiating globally for the last century or so of human rights is based on the responsibility of such geographically-defined governments to individuals within the sphere of influence of that government⁵⁵. Now wise actors like the European Union have extended the notion of their individual's sovereignty over cyber assets such as personal data⁵⁶. This makes sense for almost exactly the same reason as rights to airspace make sense. With bidirectional information access we can influence individual's behaviour just as we could with physical force.

Random House, 2017.

⁵⁴Erik Brynjolfsson, Daniel Rock, and Chad Syverson. "Artificial intelligence and the modern productivity paradox: A clash of expectations and statistics". In: *Economics of Artificial Intelligence*. University of Chicago Press, 2017.

⁵⁵Sabine C Carey, Mark Gibney, and Steven C Poe. *The politics of human rights: the quest for dignity*. Cambridge University Press, 2010.

⁵⁶Paul Nemitz. "Constitutional democracy and technology in the age of artificial intelligence". In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376.2133 (2018), p. 20180089. DOI: 10.1098/rsta.2018.0089. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/rsta.2018.0089>.

Recently there has been good reason to hope that we really will start mandating developers to follow best practice in software engineering⁵⁷. If we are sensible, we will also ensure that the information systems spreading and engulfing us will also be entirely cybersecure (or else not on the Internet), with clearly-documented accountability and lines of responsibility⁵⁸. Nevertheless, even if these visions can be achieved, there are still other areas of law and governance with which we should be concerned. The last one I focus on in this present chapter are the new centres of power and wealth. As just explained in the previous section, these are also parts of the everything human that AI and ICT are altering. Further, it is clear that achieving secure and accountable AI requires cooperation with adequate sources of power to counter those who wish to avoid the consensus of the law. Therefore wealth and power distribution, while again like cybersecurity clearly orthogonal technologically to AI, are also clearly irrevocably intertwined with its ethical and regulated application. Problems of AI accountability and grotesquely uneven wealth distribution are unlikely to be solved independently.

In this section it should be noted that I am describing my own work in progress with colleagues⁵⁹, but some of it seems sufficiently evident to justify inclusion now. For example, we hypothesise that when new technologies reduce the economic cost of distance, this in turn reduces the amount of easily-sustained competition in a sector. This is because locale becomes less a part of value, so higher quality products and services can dominate ever-larger regions, up to and including the entire globe. Such a process may have sparked the gross inequality of the late Nineteenth and Early Twentieth Centuries, when rail, news and telecommunication, and oil (far easier to transport than coal or wood) were the new monopolies. Inequality spirals if capital is allowed to capture regulation, as seems recently to have happened not only with ‘big tech’ globally, but also for example with finance in the UK or oil in Saudi Arabia and Russia, leading to a ‘resource curse’⁶⁰. In

⁵⁷OECD. *Recommendation of the Council on Artificial Intelligence*. OECD Legal Instruments OECD/LEGAL/0449. includes the OECD Principles of AI. Paris: Organisation for Economic Cooperation and Development, May 2019.

⁵⁸cf. Filippo Santoni de Sio and Jeroen van den Hoven. “Meaningful Human Control over Autonomous Systems: A Philosophical Account”. In: *Frontiers in Robotics and AI* 5 (2018), p. 15. ISSN: 2296-9144. DOI: 10.3389/frobt.2018.00015. URL: <https://www.frontiersin.org/article/10.3389/frobt.2018.00015>.

⁵⁹Alexander J Stewart, Nolan McCarty, and Joanna J Bryson. “Explaining Parochialism: A Causal Account for Political Polarization in Changing Economic Environments”. arXiv preprint arXiv:1807.11477. 2018.

⁶⁰John Christensen, Nick Shaxson, and Duncan Wigan. “The Finance Curse: Britain

the mid-Twentieth Century, stability was only reclaimed via the innovation of the welfare state, which in some countries (including the US and UK) preceded at least the second World War, though cooperation sadly there too was motivated by the first.

Governance can be almost defined by redistribution; certainly allocation of resources to solve communal problems and create public goods is governance's core characteristic⁶¹. Thus excessive inequality⁶² can be seen as a failure of governance. Right now what we are clearly not able to govern (interestingly, on both sides of the Great Firewall of China) are Internet companies. As a result, and similar to the market for commercial aircraft, the costs of distance are sufficiently negligible that the best products are very likely to become global monopolies, unless there is a substantial government investment, e.g., the Great Firewall of China⁶³, or Airbus in Europe⁶⁴. Where governance fails in a local region e.g. a county is also where we are likely to see political polarisation and populist candidates or referendum outcomes⁶⁵.

Many problems we associate with the present moment then were not necessarily created by AI or ICT directly, but rather indirectly by increasing inequality and regulatory capture. Other problems were not so much created

and the World Economy". In: *The British Journal of Politics and International Relations* 18.1 (2016), pp. 255–269. DOI: 10.1177/1369148115612793. URL: <https://doi.org/10.1177/1369148115612793>; Nolan M McCarty, Keith T Poole, and Howard Rosenthal. *Polarized America: The dance of ideology and unequal riches*. second. Cambridge, MA: MIT Press, 2016.

⁶¹Jean-Pierre Landau. "Populism and Debt: Is Europe Different from the U.S.?" Talk at the Princeton Woodrow Wilson School, and in preparation. Feb. 2016.

⁶²E.g. a Gini coefficient over 0.27 Francesco Grigoli and Adrian Robles. *Inequality Overhang*. IMF Working Paper WP/17/76. International Monetary Fund, 2017, note that too low can be problematic too.

⁶³Roya Ensafi et al. "Analyzing the Great Firewall of China over space and time". In: *Proceedings on privacy enhancing technologies* 2015.1 (2015), pp. 61–76.

⁶⁴Damien Neven and Paul Seabright. "European industrial policy: the Airbus case". In: *Economic Policy* 10.21 (July 1995), pp. 313–358. ISSN: 0266-4658. DOI: 10.2307/1344592. URL: <https://doi.org/10.2307/1344592>.

⁶⁵Yuri M. Zhukov. "Trading hard hats for combat helmets: The economics of rebellion in eastern Ukraine". In: *Journal of Comparative Economics* 44.1 (2016). Special Issue on Ukraine: Escape from Post-Soviet Legacy, pp. 1–15. ISSN: 0147-5967. DOI: <https://doi.org/10.1016/j.jce.2015.10.010>. URL: <http://www.sciencedirect.com/science/article/pii/S014759671500092X>; Sascha O Becker, Thiemo Fetzer, and Dennis Novy. "Who voted for Brexit? A comprehensive district-level analysis". In: *Economic Policy* 32.92 (Oct. 2017), pp. 601–650. ISSN: 0266-4658. DOI: 10.1093/epolic/eix012. URL: <https://doi.org/10.1093/epolic/eix012>; Florian Dorn et al. "Inequality and Extremist Voting: Evidence from Germany". In: (2018).

as exposed by AI⁶⁶. There are some exceptions where ICT—particularly, the capacity of digital media to be fully reproduced at distance inexpensively—do produce qualitative change. These include changing of the meaning of ownership⁶⁷, and generating truly novel means for recognising and disrupting human intentions, even implicit ones not known by their actors⁶⁸. On the other hand, some things are or should be treated as invariant. As an example mentioned earlier, human rights are the painstakingly agreed foundation of international law and the obligations of a state, and should be treated as core to ethical AI systems⁶⁹.

One of the disturbing things we come to understand as we learn about algorithms is the extent to which humans are ourselves algorithmic. Law can make us more so, particularly when we constrain ourselves with it, for example with mandatory sentencing. But ordinarily, humans do have wiggle room⁷⁰. As mentioned earlier, trust is based on ignorance⁷¹—that ignorance may be an important feature of society that ICT might remove. Trust allows cheating or innovating, and sometimes this may be essential. First, allowing innovation makes the level of detail about exceptions that needs to be specified more tractable. Second, of course, innovation allows us to adjust to the unexpected and to find novel, sometimes better solutions. Some—perhaps many—nations may be in danger of allowing the digital era to make innovation or free thinking too difficult or individually risky, creating nation-wide fragility to security threats as well as impinging on the important human

⁶⁶Nemitz, “Constitutional democracy and technology in the age of artificial intelligence”; Orly Mazur. “Taxing the Robots”. In: *Pepperdine Law Review* 46 (2018), pp. 277–330.

⁶⁷Aaron Perzanowski and Jason Schultz. *The End of Ownership: Personal Property in the Digital Economy*. Cambridge, MA: MIT Press, 2016.

⁶⁸Caio Machado and Marco Konopacki. “Computational Power: Automated Use of WhatsApp in the Brazilian Elections”. In: *Medium* (26 October 2018). URL: <https://feed.itsrio.org/%20computational-power-automated-use-of-whatsapp-in-the-elections-59f62b857033>; Cadwalladr, “I made Steve Bannon’s psychological warfare tool: meet the data war whistleblower”; Zhe Wu et al. “Deception detection in videos”. In: *Thirty-Second AAAI Conference on Artificial Intelligence*. 2018.

⁶⁹Philip Alston and Mary Robinson. *Human Rights and Development: Towards Mutual Reinforcement*. Oxford University Press, 2005; David Kaye. “State Execution of the International Covenant on Civil and Political Rights”. In: *UC Irvine Law Review* 3 (2013), pp. 95–125. URL: <https://scholarship.law.uci.edu/ucilr/vol13/iss1/9>.

⁷⁰Cohen, “What privacy is for”.

⁷¹O’Neill, *A question of trust: The BBC Reith Lectures 2002*; Paul Rauwolf and Joanna J. Bryson. “Expectations of Fairness and Trust Co-Evolve in Environments of Partial Information”. In: *Dynamic Games and Applications* 8.4 (Dec. 2018), pp. 891–917. ISSN: 2153-0793. DOI: 10.1007/s13235-017-0230-x. URL: <https://doi.org/10.1007/s13235-017-0230-x>.

right of freedom of opinion⁷². In such countries, law may bend too much towards the group and inadequately defend the individual. As I mentioned, this is an issue not only of rights, but also of robustness—individuals and variation produce alternatives, and choosing amongst these is a rapid way to change behaviour when a crisis demonstrates change is needed⁷³. Given that the digital revolution has fundamentally changed the nature of privacy for everyone, all societies will need to find a way to reintroduce and defend ‘wiggle room’ for innovation and opinion. I believe strongly that it would be preferable if this is done not by destroying access to history, but by acknowledging and defending individual differences, including shortcomings and the necessity of learning. But psychological and political realities remain to be explored and understood, and may vary by polity.

Summary, and the robots themselves

To reiterate my main points, when computer science is mistaken for a branch of mathematics, many important implications of computation being a physical process are lost. Further, the impact on society of the dissemination of information, power, and influence has not been adequately noted in either of those two disciplines, while in law and social sciences, awareness of technological reality and affordances has been building only slowly. Ironically, these impacts until very recently been noticed much in political science. Primarily, these impacts were most noted only in sociology, which was unfortunately in many ways imploding at the same time AI was exploding. Similarly to the myopia of computer science, psychology has primarily seen itself as studying humans as organisms, and the primary ethical considerations in that field were seen as being those of medical subjects, such as patient privacy. Again, some related disciplines such as media studies or marketing raised the issue that as we better understood human behavior we might more effectively manipulate and control it, but that observation made little headway in the popular academic understanding of artificial intelligence. Direct interventions via neuroscience and drugs received more attention, but the potential for indirect manipulations particularly of adults were seemingly dismissed.

These historic errors may be a consequence of the fact that human adults

⁷²Cf. Frischmann and Selinger, *Re-engineering humanity*, p. .

⁷³Cohen, “What privacy is for”; Luke Stark. “The emotional context of information privacy”. In: *The Information Society* 32.1 (2016), pp. 14–27. DOI: 10.1080/01972243.2015.1107167. URL: <https://doi.org/10.1080/01972243.2015.1107167>.

are of necessity the ultimate moral agents. We are the centres of accountability in our own societies, and as such are expected to have the capacity to take care of ourselves. AI ethics therefore was often reduced to its popular culture edifice as an extension of the civil rights movement⁷⁴. Now that we have discovered—astonishingly!—that people of other ethnicities and genders are as human as ‘we’ are, ‘we’ are therefore obliged by some to consider that *anything* might be human. This position seems more a rejection of the inclusivity of civil and human rights than an appropriate extension, but it is powerfully attractive to many who seem particularly likely to be members of the recently dominant gender and ethnicity, and who perhaps intuit that such an extension would again raise the power of their on clique by making the notion of rights less meaningful.

More comprehensibly, some have suggested we must extend human rights protections to anything that humans might identify with in order to protect that self concept, even if it is implicit or mistaken⁷⁵. This follows from Kant’s observation that those that treat animals reminiscent of humanity badly are also more likely to treat humans badly. Extending this principle to AI is also most likely a mistake, and an avoidable one. Remember that AI is definitionally an artefact and therefore designed. It almost certainly makes more sense where tractable to change AI than to radically change the law. The UK first⁷⁶ and now very recently the OECD⁷⁷ have recommended that AI should *not* deceptively appear to be human, so such confusions might be minimised. This may seem heavily restrictive at present, but as society becomes more familiar with AI—and through that process, better understands what it is about being human that requires and deserves protection—we should be able to broaden the scope of how near to human

⁷⁴Tony J. Prescott. “Robots are not just tools”. In: *Connection Science* 29.2 (2017), pp. 142–149. DOI: 10.1080/09540091.2017.1279125. URL: <https://doi.org/10.1080/09540091.2017.1279125>; David J Gunkel. “The other question: can and should robots have rights?” In: *Ethics and Information Technology* 20.2 (2018), pp. 87–99; Daniel Estrada. “Value Alignment, Fair Play, and the Rights of Service Robots”. In: *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. AIES ’18. New York, NY, USA: ACM, 2018, pp. 102–107. ISBN: 978-1-4503-6012-8. DOI: 10.1145/3278721.3278730. URL: <http://doi.acm.org/10.1145/3278721.3278730>.

⁷⁵Joel Parthemore and Blay Whitby. “What Makes Any Agent a Moral Agent? Reflections on Machine Consciousness and Moral Agency”. In: *International Journal of Machine Consciousness* 05.02 (2013), pp. 105–129. DOI: 10.1142/S1793843013500017. URL: <https://doi.org/10.1142/S1793843013500017>; David J Gunkel. *Robot rights*. MIT Press, 2018.

⁷⁶Boden et al., *Principles of Robotics*.

⁷⁷OECD, *Recommendation of the Council on Artificial Intelligence*.

devices can be while still not having them be deceptive⁷⁸.

As discussed earlier, there are recent calls to ground AI not on ‘ethics’ (which is viewed as ill defined) but on international human rights law. Of course, this may be a false dichotomy; procedures from classical ethics theories may still be of use in determining ambiguities and tradeoffs of application⁷⁹. We can certainly expect ongoing consideration of localised variation which *ethics* perhaps better encapsulates than *rights*. Ethics has always been about codes of conduct, which confound fundamental principles which we may be able to codify with rights, with other things that are essentially identity markers. But identity too can be essential to security through constructing a defensible identity⁸⁰. Identity obviously (definitionally) defines a group, and groups are often the best means humans have for achieving security and therefore viability. Not only is breaking into different groups sometimes more efficient for governance or other resource constraints, but also some groups will have different fundamental security tradeoffs based on their geological and ecological situation, and also just simply their neighbours. That identity often also rests on shared historical narratives which will afford different organisational strategies as well may be secondary to the more essential geo-ecological concerns (as is illustrated by the apparent ease with which new ethnicities are invented⁸¹) it still of course also makes a contribution.

In conclusion, any artefact that transforms perception to more relevant information, including action, is AI—and note that AI is an adjective, not a noun, unless it is referring to the academic discipline. There is no question that AI and digital technologies more generally are introducing enormous transformations to society. AI Nevertheless, these impacts should be governable by less transformative legislative change. The vast majority of AI—

⁷⁸Joanna J. Bryson. “The Meaning of the EPSRC Principles of Robotics”. In: *Connection Science* 29.2 (2017), pp. 130–136. DOI: 10.1080/09540091.2017.1313817. URL: <http://dx.doi.org/10.1080/09540091.2017.1313817>.

⁷⁹Cansu Canca. “Human Rights and AI Ethics: Why Ethics Cannot be Replaced by the UDHR”. in: *United Nations University: AI & Global Governance Articles & Insights* (July 2019). URL: <https://cpr.unu.edu/%20ai-global-governance-human-rights-and-ai-ethics-why-ethics-cannot-be-replaced-by-the-udhr.html>.

⁸⁰Bill McSweeney. *Security, identity and interests: a sociology of international relations*; Simon T. Powers. “The Institutional Approach for Modeling the Evolution of Human Societies”. In: *Artificial Life* 24.1 (2018). PMID: 29369715, pp. 10–28. DOI: 10.1162/ARTL_a_00251. URL: https://doi.org/10.1162/ARTL_a_00251.

⁸¹Erin K. Jenne, Stephen M. Saideman, and Will Lowe. “Separatism as a Bargaining Posture: The Role of Leverage in Minority Radicalization”. In: *Journal of Peace Research* 44.5 (2007), pp. 539–558. DOI: 10.1177/0022343307080853. URL: <https://doi.org/10.1177/0022343307080853>.

particularly where it has social impact is and will remain a consequence of corporate commercial processes, and as such subject to existing regulations and regulating strategies. We may require more regulatory bodies with expertise in examining the accounts of software development, but it is critical to remember that what we are holding accountable is not machines themselves, but the people who build, own, or operate them—including any who alter their operation through assault on their cybersecurity. What we need to govern is the human application of technology; what we need to oversee are human processes of development, testing, operation, and monitoring.

AI also offered us an opportunity to discover more about how we ourselves and our societies work. By allowing us to construct artefacts that mimic aspects of Nature but with new affordances for modularity and decoupling, we allow ourselves novel means of self examination, including of our most crucial capacities such as morality and political behaviour. This is an exciting time for scientific and artistic exploration as well as for commerce and law. But better knowledge also offers an opportunity for better control. The role of the law for crafting both individual and societal protections has never been more crucial.

Acknowledgements

A small proportion of the material in this review was derived from a document previously delivered to the OECD (Karine Perset) in May 2017 under the title “Current and Potential Impacts of Artificial Intelligence and Autonomous Systems on Society,” which contributed to the OECD AI policy efforts and documents of 2018–2019, and also reused (with permission) and expanded for. More debt is probably owed to Frank Pasquale for extremely useful feedback and suggestions on a first draft. Thanks also to Will Lowe, Patrick Slavenburg, and Jean-Paul Skeete. I was supported in part by an AXA Research Fellowship in AI Ethics while writing this chapter.

References

- Alston, Philip and Mary Robinson. *Human Rights and Development: Towards Mutual Reinforcement*. Oxford University Press, 2005.
- Andrejevic, Mark. “Automating Surveillance”. In: *Surveillance & Society* 17.1/2 (2019), pp. 7–13.

⁸¹Bryson, “The Past Decade and Future of AI’s Impact on Society”.

- Armstrong, Harvey and Robert Read. “Western European micro-states and EU autonomous regions: The advantages of size and sovereignty”. In: *World Development* 23.7 (1995), pp. 1229–1245. ISSN: 0305-750X. DOI: [http://dx.doi.org/10.1016/0305-750X\(95\)00040-J](http://dx.doi.org/10.1016/0305-750X(95)00040-J). URL: <http://www.sciencedirect.com/science/article/pii/S0305750X9500040J>.
- Barr, Joseph R. and Joseph Cavanaugh. “Forensics: Assessing model goodness: A machine learning view”. In: *Robotic Intelligence*. 2019, pp. 17–23. DOI: 10.1142/9789811203480_0003. URL: https://www.worldscientific.com/doi/abs/10.1142/9789811203480_0003.
- Becker, Sascha O, Thiemo Fetzer, and Dennis Novy. “Who voted for Brexit? A comprehensive district-level analysis”. In: *Economic Policy* 32.92 (Oct. 2017), pp. 601–650. ISSN: 0266-4658. DOI: 10.1093/epolic/eix012. URL: <https://doi.org/10.1093/epolic/eix012>.
- Bertrand, Marianne and Sendhil Mullainathan. “Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination”. In: *The American Economic Review* 94.4 (2004), pp. 991–1013.
- Bloom, Paul. *Against empathy: The case for rational compassion*. Random House, 2017.
- Boden, Margaret A. “Robot says: Whatever (The robots won’t take over because they couldn’t care less)”. In: *Aeon* (13 August 2018). originally a lecture at the Leerhulme Centre for the Future of Intelligence. URL: <https://aeon.co/essays/the-robots-wont-take-over-because-they-couldnt-care-less>.
- Boden, Margaret et al. *Principles of Robotics*. The United Kingdom’s Engineering and Physical Sciences Research Council (EPSRC). Apr. 2011. URL: <https://www.epsrc.ac.uk/research/ourportfolio/themes/engineering/activities/principlesofrobotics/>.
- Bostrom, Nick. *Superintelligence: Paths, dangers, strategies*. Oxford University Press, 2014.
- Box, G. E. P. “Robustness in the strategy of scientific model building”. In: *Robustness in statistics*. Ed. by R. L. Launer and G. N. Wilkinson. New York, NY: Academic Press, 1979, pp. 201–236.
- Brundage, Miles and Joanna J. Bryson. *Smart Policies for Artificial Intelligence*. in preparation, available as arXiv:1608.08196. 2017.
- Brundage, Miles et al. *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation*. Tech. rep. <https://maliciousaireport.com/>. Future of Humanity Institute, University of Oxford, Centre for the Study of Existential Risk, University of Cambridge, Center for a New American Security, Electronic Frontier Foundation, and OpenAI, Feb. 2018.

- Brynjolfsson, Erik, Daniel Rock, and Chad Syverson. “Artificial intelligence and the modern productivity paradox: A clash of expectations and statistics”. In: *Economics of Artificial Intelligence*. University of Chicago Press, 2017.
- Bryson, Joanna J. “Embodiment versus Memetics”. In: *Mind & Society* 7.1 (June 2008), pp. 77–94.
- “Artificial Intelligence and Pro-Social Behaviour”. In: *Collective Agency and Cooperation in Natural and Artificial Systems: Explanation, Implementation and Simulation*. Ed. by Catrin Misselhorn. Vol. 122. Philosophical Studies. Berlin: Springer, Oct. 2015, pp. 281–306.
- “The Meaning of the EPSRC Principles of Robotics”. In: *Connection Science* 29.2 (2017), pp. 130–136. DOI: 10.1080/09540091.2017.1313817. URL: <http://dx.doi.org/10.1080/09540091.2017.1313817>.
- “Patiency is not a virtue: the design of intelligent systems and systems of ethics”. In: *Ethics and Information Technology* 20.1 (Mar. 2018), pp. 15–26. ISSN: 1572-8439. DOI: 10.1007/s10676-018-9448-6. URL: <https://doi.org/10.1007/s10676-018-9448-6>.
- “The Past Decade and Future of AI’s Impact on Society”. In: *Towards a New Enlightenment? A Transcendent Decade*. OpenMind BBVA. commissioned, based on a white paper also commissioned, that by the OECD. Madrid: Taylor, Mar. 2019. URL: <https://www.bbvaopenmind.com/en/articles/%20the-past-decade-and-future-of-ais-impact-on-society/>.
- Bryson, Joanna J., Mihailis E. Diamantis, and Thomas D. Grant. “Of, for, and by the people: the legal lacuna of synthetic persons”. In: *Artificial Intelligence and Law* 25.3 (Sept. 2017), pp. 273–291. ISSN: 1572-8382. DOI: 10.1007/s10506-017-9214-9. URL: <https://doi.org/10.1007/s10506-017-9214-9>.
- Bryson, Joanna J. and Alan F. T. Winfield. “Standardizing Ethical Design for Artificial Intelligence and Autonomous Systems”. In: *Computer* 50.5 (May 2017), pp. 116–119. ISSN: 0018-9162. DOI: 10.1109/MC.2017.154.
- Čače, Ivana and Joanna J. Bryson. “Agent Based Modelling of Communication Costs: Why Information can be Free”. In: *Emergence and Evolution of Linguistic Communication*. Ed. by C. Lyon, C. L. Nehaniv, and A. Cangelosi. London: Springer, 2007, pp. 305–322.
- Cadwalladr, Carole. “‘I made Steve Bannon’s psychological warfare tool’: meet the data war whistleblower”. In: *The Observer* (18 March 2018).
- Calinon, Sylvain et al. “Learning and reproduction of gestures by imitation”. In: *IEEE Robotics & Automation Magazine* 17.2 (2010), pp. 44–54.

- Caliskan, Aylin, Joanna J. Bryson, and Arvind Narayanan. “Semantics derived automatically from language corpora contain human-like biases”. In: *Science* 356.6334 (2017), pp. 183–186. ISSN: 0036-8075. DOI: 10.1126/science.aal4230. URL: <http://science.sciencemag.org/content/356/6334/183>.
- Canca, Cansu. “Human Rights and AI Ethics: Why Ethics Cannot be Replaced by the UDHR”. In: *United Nations University: AI & Global Governance Articles & Insights* (July 2019). URL: <https://cpr.unu.edu/%20ai-global-governance-human-rights-and-ai-ethics-why-ethics-cannot-be-replaced-by-the-udhr.html>.
- Carey, Sabine C, Mark Gibney, and Steven C Poe. *The politics of human rights: the quest for dignity*. Cambridge University Press, 2010.
- Choe, Yoonsuck, Jaerock Kwon, and Ji Ryang Chung. “Time, Consciousness, and Mind Uploading”. In: *International Journal of Machine Consciousness* 04.01 (2012), pp. 257–274. DOI: 10.1142/S179384301240015X. URL: <http://www.worldscientific.com/doi/abs/10.1142/%20S179384301240015X>.
- Christensen, John, Nick Shaxson, and Duncan Wigan. “The Finance Curse: Britain and the World Economy”. In: *The British Journal of Politics and International Relations* 18.1 (2016), pp. 255–269. DOI: 10.1177/1369148115612793. URL: <https://doi.org/10.1177/1369148115612793>.
- Claxton, Guy. *Intelligence in the flesh: Why your mind needs your body much more than it thinks*. Yale University Press, 2015.
- Cohen, Julie E. “What privacy is for”. In: *Harvard Law Review* 126 (May 2013), pp. 1904–1933.
- Cooke, Maeve. “A space of one’s own: Autonomy, privacy, liberty”. In: *Philosophy & Social Criticism* 25.1 (1999), pp. 22–53. DOI: 10.1177/019145379902500102. URL: <http://dx.doi.org/10.1177/019145379902500102>.
- Cranefield, Stephen et al. “No Pizza for You: Value-based Plan Selection in BDI Agents”. In: *IJCAI Proceedings*. Ed. by Carles Sierra. Melbourne, Aug. 2017.
- Dennett, Daniel C. “Why You Can’t Make a Computer that Feels Pain”. In: *Brainstorms*. page numbers are from the 1986 Harvester Press Edition, Brighton, Sussex. Montgomery, Vermont: Bradford Books, 1978, pp. 190–229.
- *From Bacteria to Bach and Back*. Allen Lane, 2017.
- Dorn, Florian et al. “Inequality and Extremist Voting: Evidence from Germany”. In: (2018).
- Dredze, Mark et al. “Intelligent email: Reply and attachment prediction”. In: *Proceedings of the 13th international conference on Intelligent user interfaces*. ACM. 2008, pp. 321–324.

- Ensafi, Roya et al. “Analyzing the Great Firewall of China over space and time”. In: *Proceedings on privacy enhancing technologies* 2015.1 (2015), pp. 61–76.
- Erdem. “The trade-off in machine learning: Accuracy vs explainability”. In: *Medium* (Dec. 2018). URL: <https://medium.com/@erdemkalayci/%20the-tradeoff-in-machine-learning-accuracy-vs-explainability-fbb13914fde2>.
- Estrada, Daniel. “Value Alignment, Fair Play, and the Rights of Service Robots”. In: *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. AIES ’18. New York, NY, USA: ACM, 2018, pp. 102–107. ISBN: 978-1-4503-6012-8. DOI: 10.1145/3278721.3278730. URL: <http://doi.acm.org/10.1145/3278721.3278730>.
- Frischmann, Brett and Evan Selinger. *Re-engineering humanity*. Cambridge University Press, 2018.
- Grigoli, Francesco and Adrian Robles. *Inequality Overhang*. IMF Working Paper WP/17/76. International Monetary Fund, 2017.
- Gunkel, David J. *Robot rights*. MIT Press, 2018.
- “The other question: can and should robots have rights?” In: *Ethics and Information Technology* 20.2 (2018), pp. 87–99.
- Hall, Patrick. “On the Art and Science of Machine Learning Explanations”. In: *arXiv preprint arXiv:1810.02909* (2018).
- Hamilton, Barton H., Jack A. Nickerson, and Hideo Owan. “Diversity and Productivity in Production Teams”. In: *Advances in the Economic Analysis of Participatory and Labor-Managed Firms*. 2012, pp. 99–138. DOI: 10.1108/S0885-3339(2012)0000013009. URL: <https://www.emeraldinsight.com/doi/abs/10.1108/%20S0885-3339%282012%290000013009>.
- Hoven, Jeroen van den. “ICT and Value Sensitive Design”. In: *The Information Society: Innovation, Legitimacy, Ethics and Democracy In honor of Professor Jacques Berleur s.j.* Ed. by Philippe Goujon et al. Boston, MA: Springer US, 2007, pp. 67–72. ISBN: 978-0-387-72381-5.
- Hüttermann, Michael. *DevOps for Developers*. Apress, Springer, 2012.
- Jenne, Erin K., Stephen M. Saideman, and Will Lowe. “Separatism as a Bargaining Posture: The Role of Leverage in Minority Radicalization”. In: *Journal of Peace Research* 44.5 (2007), pp. 539–558. DOI: 10.1177/0022343307080853. URL: <https://doi.org/10.1177/0022343307080853>.
- Jordan, Jillian J. et al. “Uncalculating cooperation is used to signal trustworthiness”. In: *Proceedings of the National Academy of Sciences* (2016). DOI: 10.1073/pnas.1601280113. URL: <http://www.pnas.org/content/early/2016/07/19/%201601280113.abstract>.

- Katzman, Kenneth and Paul K Kerr. *Iran nuclear agreement*. Tech. rep. R43333. Library of Congress, Congressional Research Service, May 2016. URL: www.crs.gov.
- Kaye, David. “State Execution of the International Covenant on Civil and Political Rights”. In: *UC Irvine Law Review* 3 (2013), pp. 95–125. URL: <https://scholarship.law.uci.edu/ucilr/vol3/iss1/9>.
- Kozyrkov, Cassie. “Explainable AI won’t deliver. Here’s why.” In: *Hackernoon* (Nov. 2018). URL: <https://hackernoon.com/%20explainable-ai-wont-deliver-here-s-why-6738f54216be>.
- Kroll, Joshua A. et al. “Accountable Algorithms”. In: *University of Pennsylvania Law Review* 165 (2017), pp. 633–706.
- Kshemkalyani, Ajay D and Mukesh Singhal. *Distributed computing: principles, algorithms, and systems*. Cambridge University Press, 2011.
- Landau, Jean-Pierre. “Populism and Debt: Is Europe Different from the U.S.?” Talk at the Princeton Woodrow Wilson School, and in preparation. Feb. 2016.
- Lewis, Molly and Gary Lupyan. “Language use shapes cultural norms: Large scale evidence from gender”. In: *Proceedings of the Annual Meeting of the Cognitive Science Society (CogSci)*. also in prep. for journal publication. Madison, WI, 2018, pp. 2041–2046.
- List, Christian and Philip Pettit. *Group agency: The possibility, design, and status of corporate agents*. Oxford University Press, 2011.
- Mace, Ruth. “The co-evolution of human fertility and wealth inheritance strategies”. In: *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 353.1367 (1998), pp. 389–397. ISSN: 0962-8436. DOI: 10.1098/rstb.1998.0217. URL: <http://rstb.royalsocietypublishing.org/content/353/1367/%20389>.
- Machado, Caio and Marco Konopacki. “Computational Power: Automated Use of WhatsApp in the Brazilian Elections”. In: *Medium* (26 October 2018). URL: <https://feed.itsrio.org/%20computational-power-automated-use-of-whatsapp-in-the-elections-59f62b857033>.
- Mazur, Orly. “Taxing the Robots”. In: *Pepperdine Law Review* 46 (2018), pp. 277–330.
- McCarty, Nolan M, Keith T Poole, and Howard Rosenthal. *Polarized America: The dance of ideology and unequal riches*. second. Cambridge, MA: MIT Press, 2016.
- McSweeney, Bill. *Security, identity and interests: a sociology of international relations*.

- Mesoudi, Alex, Andrew Whiten, and Kevin N. Laland. “Towards a unified science of cultural evolution”. In: *Behavioral and Brain Sciences* 29.4 (2006), pp. 329–347. DOI: 10.1017/S0140525X06009083.
- Moeslund, Thomas B and Erik Granum. “A survey of computer vision-based human motion capture”. In: *Computer vision and image understanding* 81.3 (2001), pp. 231–268.
- Nemitz, Paul. “Constitutional democracy and technology in the age of artificial intelligence”. In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376.2133 (2018), p. 20180089. DOI: 10.1098/rsta.2018.0089. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/rsta.2018.0089>.
- Neven, Damien and Paul Seabright. “European industrial policy: the Airbus case”. In: *Economic Policy* 10.21 (July 1995), pp. 313–358. ISSN: 0266-4658. DOI: 10.2307/1344592. URL: <https://doi.org/10.2307/1344592>.
- O’Neil, Cathy. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Broadway Books, 2016.
- O’Neill, Onora. *A question of trust: The BBC Reith Lectures 2002*. Cambridge University Press, 2002.
- O’Reilly, Tim. *WTF? What’s the Future and why It’s Up to Us*. New York: Random House, 2017.
- OECD. *Recommendation of the Council on Artificial Intelligence*. OECD Legal Instruments OECD/LEGAL/0449. includes the OECD Principles of AI. Paris: Organisation for Economic Cooperation and Development, May 2019.
- Parthemore, Joel and Blay Whitby. “What Makes Any Agent a Moral Agent? Reflections on Machine Consciousness and Moral Agency”. In: *International Journal of Machine Consciousness* 05.02 (2013), pp. 105–129. DOI: 10.1142/S1793843013500017. URL: <https://doi.org/10.1142/S1793843013500017>.
- Pasquale, Frank. “Two concepts of immortality: Reframing public debate on stem-cell research”. In: *Yale Journal of Law and the Humanities* 14 (2002), pp. 73–121.
- Perzanowski, Aaron and Jason Schultz. *The End of Ownership: Personal Property in the Digital Economy*. Cambridge, MA: MIT Press, 2016.
- Powers, Simon T. “The Institutional Approach for Modeling the Evolution of Human Societies”. In: *Artificial Life* 24.1 (2018). PMID: 29369715, pp. 10–28. DOI: 10.1162/ARTL_a_00251. URL: https://doi.org/10.1162/ARTL_a_00251.

- Powers, Simon T., Carel P. van Schaik, and Laurent Lehmann. “How institutions shaped the last major evolutionary transition to large-scale human societies”. In: *Philosophical Transactions of the Royal Society B: Biological Sciences* 371.1687 (2016), p. 20150098. DOI: 10.1098/rstb.2015.0098. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.2015.0098>.
- Prescott, Tony J. “Robots are not just tools”. In: *Connection Science* 29.2 (2017), pp. 142–149. DOI: 10.1080/09540091.2017.1279125. URL: <https://doi.org/10.1080/09540091.2017.1279125>.
- Rauwolf, Paul and Joanna J. Bryson. “Expectations of Fairness and Trust Co-Evolve in Environments of Partial Information”. In: *Dynamic Games and Applications* 8.4 (Dec. 2018), pp. 891–917. ISSN: 2153-0793. DOI: 10.1007/s13235-017-0230-x. URL: <https://doi.org/10.1007/s13235-017-0230-x>.
- Romanes, George John. *Animal intelligence*. London: D. Appleton, 1882.
- Roughgarden, Joan, Meeko Oishi, and Erol Akçay. “Reproductive Social Behavior: Cooperative Games to Replace Sexual Selection”. In: *Science* 311.5763 (2006), pp. 965–969. DOI: 10.1126/science.1110105. URL: <http://www.sciencemag.org/content/311/5763/965.abstract>.
- Santoni de Sio, Filippo and Jeroen van den Hoven. “Meaningful Human Control over Autonomous Systems: A Philosophical Account”. In: *Frontiers in Robotics and AI* 5 (2018), p. 15. ISSN: 2296-9144. DOI: 10.3389/frobt.2018.00015. URL: <https://www.frontiersin.org/article/10.3389/frobt.2018.00015>.
- Shanahan, Murray. *The technological singularity*. MIT Press, 2015.
- Shannon, Claude Elwood. “A mathematical theory of communication”. In: *Bell system technical journal* 27.3 (1948), pp. 379–423.
- Shi, Feng et al. “The wisdom of polarized crowds”. In: *Nature Human Behaviour* 3 (2019), pp. 329–336.
- Sipser, Michael. *Introduction to the Theory of Computation*. Second. Boston, MA: PWS, Thompson, 2005.
- Smith, Kenny and Elizabeth Wonnacott. “Eliminating unpredictable variation through iterated learning”. In: *Cognition* 116.3 (2010), pp. 444–449. ISSN: 0010-0277. DOI: <https://doi.org/10.1016/j.cognition.2010.06.004>. URL: <http://www.sciencedirect.com/science/article/pii/S0010027710001320>.
- Stark, Luke. “The emotional context of information privacy”. In: *The Information Society* 32.1 (2016), pp. 14–27. DOI: 10.1080/01972243.2015.1107167. URL: <https://doi.org/10.1080/01972243.2015.1107167>.

- Stewart, Alexander J, Nolan McCarty, and Joanna J Bryson. “Explaining Parochialism: A Causal Account for Political Polarization in Changing Economic Environments”. arXiv preprint arXiv:1807.11477. 2018.
- Stoddart, David Michael. *The Scented Ape: The Biology and Culture of Human Odour*. Cambridge University Press, Nov. 1990.
- Traum, David et al. “New Dimensions in Testimony: Digitally preserving a Holocaust survivor’s interactive storytelling”. In: *Proceedings of the Eighth International Conference on Interactive Digital Storytelling*, pp. 269–281.
- Wemelsfelder, Françoise. “The scientific validity of subjective concepts in models of animal welfare”. In: *Applied Animal Behaviour Science* 53.1 (1997). Basic and Applied Aspects of Motivation and Cognition, pp. 75–88. ISSN: 0168-1591. DOI: [https://doi.org/10.1016/S0168-1591\(96\)01152-5](https://doi.org/10.1016/S0168-1591(96)01152-5). URL: <http://www.sciencedirect.com/science/article/pii/S0168159196011525>.
- Williams Woolley, Anita et al. “Evidence for a Collective Intelligence Factor in the Performance of Human Groups”. In: *Science* 330.6004 (29 October 2010), pp. 686–688.
- Wu, Zhe et al. “Deception detection in videos”. In: *Thirty-Second AAAI Conference on Artificial Intelligence*. 2018.
- Wynsberghe, Aimee van. “Designing Robots for Care: Care Centered Value-Sensitive Design”. In: *Science and Engineering Ethics* 19.2 (June 2013), pp. 407–433. ISSN: 1471-5546. DOI: [10.1007/s11948-011-9343-6](https://doi.org/10.1007/s11948-011-9343-6). URL: <https://doi.org/10.1007/s11948-011-9343-6>.
- Zeng, Jiaming, Berk Ustun, and Cynthia Rudin. “Interpretable classification models for recidivism prediction”. In: *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 180.3 (2017), pp. 689–722. ISSN: 1467-985X. DOI: [10.1111/rssa.12227](https://doi.org/10.1111/rssa.12227). URL: <http://dx.doi.org/10.1111/rssa.12227>.
- Zhukov, Yuri M. “Trading hard hats for combat helmets: The economics of rebellion in eastern Ukraine”. In: *Journal of Comparative Economics* 44.1 (2016). Special Issue on Ukraine: Escape from Post-Soviet Legacy, pp. 1–15. ISSN: 0147-5967. DOI: <https://doi.org/10.1016/j.jce.2015.10.010>. URL: <http://www.sciencedirect.com/science/article/pii/S014759671500092X>.