

*Annual Review of Psychology*

The Moral Psychology of Artificial Intelligence

Jean-François Bonnefon,<sup>1</sup> Iyad Rahwan,<sup>2</sup> and Azim Shariff<sup>3</sup>

<sup>1</sup>Centre National de la Recherche Scientifique (TSM-R), Toulouse School of Economics, Toulouse, France; email: jean-francois.bonnefon@tse-fr.eu

<sup>2</sup>Center for Humans and Machines, Max Planck Institute for Human Development, Berlin, Germany

<sup>3</sup>Department of Psychology, University of British Columbia, Vancouver, British Columbia, Canada

ANNUAL REVIEWS **CONNECT**

[www.annualreviews.org](http://www.annualreviews.org)

- Download figures
- Navigate cited references
- Keyword search
- Explore related articles
- Share via email or social media

Annu. Rev. Psychol. 2024. 75:653–75

First published as a Review in Advance on September 18, 2023

The *Annual Review of Psychology* is online at [psych.annualreviews.org](http://psych.annualreviews.org)

<https://doi.org/10.1146/annurev-psych-030123-113559>

Copyright © 2024 by the author(s). This work is licensed under a Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited. See credit lines of images or other third-party material in this article for license information.

**Keywords**

psychology, morality, artificial intelligence, AI, agent, patient, algorithmic bias, blame, ethical dilemmas, value alignment, cooperation, humanization, delegation, AI-mediated communication

**Abstract**

Moral psychology was shaped around three categories of agents and patients: humans, other animals, and supernatural beings. Rapid progress in artificial intelligence has introduced a fourth category for our moral psychology to deal with: intelligent machines. Machines can perform as moral agents, making decisions that affect the outcomes of human patients or solving moral dilemmas without human supervision. Machines can be perceived as moral patients, whose outcomes can be affected by human decisions, with important consequences for human–machine cooperation. Machines can be moral proxies that human agents and patients send as their delegates to moral interactions or use as a disguise in these interactions. Here we review the experimental literature on machines as moral agents, moral patients, and moral proxies, with a focus on recent findings and the open questions that they suggest.

## Contents

1. INTRODUCTION .....	654
2. MACHINES AS MORAL AGENTS .....	655
2.1. Implicit Moral Machines .....	655
2.2. Explicit Moral Machines .....	658
3. MACHINES AS MORAL PATIENTS .....	662
3.1. Machine-Regarding Preferences .....	663
3.2. Overcoming the Machine Penalty .....	665
4. MACHINES AS MORAL PROXIES .....	666
4.1. Delegation to Machines .....	666
4.2. Machine Masquerade .....	667
5. CONCLUSION .....	668

## 1. INTRODUCTION

Encounters between humans and artificial intelligence (AI) have, until recently, been confined to science fiction. Droids and replicants, Commander Data and Agent Smith, T-800 and HAL 9000 have all prodded people to consider the moral questions that arise when people interact with advanced machines capable of human-level intelligence. How ought they be treated? How will they treat us? And how do they change how we treat each other?

Today, AI has finally moved beyond fiction and begun its march toward ubiquity. With the pace of AI innovations beginning to be measured in months rather than years, there is a feeling of being in the foothills of the long-promised AI revolution. As of writing, generative AI and large language models have redoubled public interest in AI. Virtual assistants are everyday tools. Recommendation algorithms govern our attention.

As human encounters with robots and other AI have multiplied, so have efforts to understand the moral psychology of AI. Some of this work involves speculative research about the not-yet-possible. Decision making on how driverless cars ought to be programmed to distribute risk to different individuals is one widely discussed example. What has been more common, however, is a process of catch-up in which researchers have been racing to understand the psychological dimensions that accompany the rapidly emerging innovations in the AI space. Online bots, AI-assisted medical diagnoses, and the use of predictive algorithms for policing and incarceration have all provoked important moral questions about trust, bias, and value alignment.

In this article, we review the research on the moral roles that intelligent machines have begun to occupy. As moral agents (see the sidebar titled *The Moral Typcasting of Machines*), machines are implicitly or explicitly charged with contributing to or making moral decisions—often about matters of life and death such as who deserves a kidney transplant, whose safety should be prioritized in traffic collisions, or who goes to jail. How should we align AI-driven decisions in these domains with human values? As moral patients, machines are the subjects of human moral behavior, be it cooperative or competitive, sympathetic or malicious. Although considering the patiency of nonsentient machines may sound like fanciful flirtation with science fiction, figuring out how to increase human cooperation with AI is already a present-day challenge. Finally, in the role of moral proxies, machines may serve as moral intermediaries in people’s treatment of their fellow humans. In this role, people can use machines to disguise, whitewash, or carry out their morally questionable behavior.

## THE MORAL TYPECASTING OF MACHINES

In this article, we speak of machines as moral agents or patients purely for the purpose of organizing the empirical findings of moral psychology; we do not mean that experts should use these terms the way we do, or that laypersons do use these terms the way we do. Our use of these terms does not imply any ontological commitment; we have nothing to contribute to philosophical debates about whether it is appropriate to attribute agency or patiency to a machine. Furthermore, our use of these terms does not imply that people engage in a binary classification of machines as agents or patients. Agency and patiency are two continuous dimensions of mind perception (Wegner & Gray 2016), and people can perceive machines as occupying various positions along that two-dimensional space.

## 2. MACHINES AS MORAL AGENTS

### 2.1. Implicit Moral Machines

A machine can be perceived as a moral agent even when its programming does not explicitly encode moral values, as long as the consequences of its actions can fall in the moral domain. This is what we call an implicit moral agent (Moor 2006) or an implicit moral machine. The prototypical case here is that of a machine whose mistakes can create harm. For example, medical AI can harm patients by making a wrong diagnosis, a recommendation algorithm can create harm by steering a child to a violent video, and a face recognition algorithm can harm a person by mistaking them for a known terrorist. These implicit moral machines are not necessarily trying to solve moral dilemmas, but their failures have moral implications. Accordingly, from a moral psychology perspective, their performance is the most important consideration. Here we consider in turn the expectations that people have about the performance of machines whose mistakes can create harm and their reactions to these mistakes. More specifically, we consider the number of mistakes people are willing to tolerate from implicit moral machines, their concerns about the distribution of these mistakes across vulnerable and less vulnerable groups, and the blame they direct toward machines that fail alone or in conjunction with humans.

**2.1.1. Performance.** How many crashes are you willing to tolerate from self-driving cars, per million kilometers? How many mistakes are you willing to accept from a skin cancer detection algorithm, per million patients? These are very hard questions. An easy way out would be to answer “zero,” that is, to require perfect performance from a machine before it is allowed to replace humans—but this extreme position would forfeit the benefits that machines can deliver even before they are fail-proof. If we require self-driving cars to be perfectly safe before we allow them on the road, we sacrifice the thousands of lives they could have saved by being allowed on the road just a little sooner (Kalra & Groves 2017). If we wait for skin cancer detection algorithms to be perfectly accurate, we sacrifice the thousands of lives that could have been saved by an earlier detection (Esteva et al. 2017). As a result, we may have a moral imperative to allow machines to make some mistakes and a need to decide how many we will allow.

Generally speaking, we may not want to let machines make decisions if they make more harmful mistakes than humans do—and, conversely, we may be willing to let machines make decisions as soon as they make fewer harmful mistakes than humans do. This is the approach taken in several policy reports about autonomous driving, which suggest that the minimal requirement before deploying self-driving cars on our roads is that they are provably safer than the average human driver (Bonneton et al. 2020a, Luetge 2017, Santoni de Sio 2021). Providing objective evidence that a machine performs better than humans, however, is not trivial to begin with (Kalra

& Paddock 2016, Kleinberg et al. 2018, Noy et al. 2018), and psychological biases may complicate things even further.

Indeed, it appears that people have extreme performance requirements for implicit moral machines, because they expect a substantial increase over baseline human performance, while overestimating this baseline human performance. For example, a representative sample of the German population believed that human experts would have a 20–30% mistake rate when predicting credit default or recidivism, which is probably an underestimation—and working from this baseline, the same sample required machines to have a mistake rate lower than 10% (Rebitschek et al. 2021). An even stronger bias exists in the domain of autonomous driving (Liu et al. 2019b, Shariff et al. 2021), where people require their self-driving car to be significantly safer than they themselves are, while substantially overestimating the safety of their own driving. In a representative sample of US drivers, the median respondent believed themselves to be in the top 25% of drivers and estimated that two-thirds of car crashes would be avoided if everyone drove like them. From this baseline, they required very high safety from self-driving cars, way above the actual average safety of human drivers.

Similar findings are available for other, less quantifiable aspects of human performance. For example, one of the main concerns that Americans have about implicit moral machines is that they do not understand nuance and complexity as well as humans do (Smith 2019). This concern translates into resistance to medical AI: Because patients think their unique characteristics and circumstances will be poorly understood by AI, they prefer to turn to human doctors (Longoni et al. 2019), leaving unexamined the actual ability of human doctors to take into account these unique characteristics and circumstances. In like vein, people express concerns about the transparency or intelligibility of medical AI recommendations, compared to that of human doctors, without realizing that they overestimate their ability to understand human doctors in the first place (Cadario et al. 2021).

**2.1.2. Bias.** Because implicit moral machines can harm people through their mistakes, people are rightly concerned about how many mistakes they make. The distribution of these mistakes also matters. Beyond how many mistakes they make, it matters whether a credit-scoring algorithm makes more mistakes about women than about men (Bono et al. 2021, Hassani 2021); it matters whether self-driving cars are less likely to detect and protect pedestrians than other road users (Combs et al. 2019); and it matters whether face recognition algorithms are more likely to misclassify dark-skinned faces (Birhane 2022, Buolamwini & Gebre 2018). The nature of the mistakes matters, too. In a landmark investigation (Angwin et al. 2016), the news organization *ProPublica* published evidence of a racial bias in the results of the COMPAS algorithm, which is used in some US courts to predict (among other outcomes) the risk that a defendant will be re-arrested in the next two years. The key result of the analysis was that while the algorithm made the same number of mistakes for black defendants and for white defendants, it did not make the same mistakes: Mistakes that were favorable to the defendant were more likely when the defendant was white, and mistakes that were unfavorable to the defendant were more likely when the defendant was black. Comparable results were later found for white versus Hispanic defendants (Hamilton 2019).

This analysis was probably the catalyst for a surge of interest in the design of algorithms whose outcomes satisfy some mathematical definition of fairness across individuals and groups. Much of this literature on algorithmic fairness is grounded in computer science and impossibility theorems, dealing with the problem that there are many possible mathematical definitions of fairness, whose requirements are sometimes impossible to achieve simultaneously (for entry points, see Chouldechova 2017, Kleinberg et al. 2017, Pleiss et al. 2017). Given that not all forms of fairness

are simultaneously achievable, it may seem natural to collect experimental data on the forms of fairness that people prefer. This experimental work is mostly disconnected from moral psychology (for a review, see Starke et al. 2022), and its results seem to be highly dependent on the application domain considered in each article. For example, people seem to prefer simple demographic parity when considering university admission algorithms, that is, to require similar admission rates for all demographic groups of applicants (Srivastava et al. 2019); when algorithms decide whether to grant bail to defendants, people prefer that they equalize false positive rates across groups rather than accuracy across groups (Harrison et al. 2020); and when algorithms decide how to allocate loans, people prefer that they adopt some calibrated version of fairness that prioritizes applicants with the highest payback rates (Saxena et al. 2020).

In view of this variation in findings across experimental protocols and domains of application, there seem to be great opportunities for designing methodologically systematic, psychology-driven programs about the kind of fairness people want from machines whose decisions can have a disparate impact across groups. In parallel, research is needed to better understand the concerns that people have about algorithmic fairness. At first sight, there are plenty of reasons to expect people to feel deep concern. First, there is ample discussion in the media about the danger that machines will learn, amplify, and legitimize the biases embedded in the human decisions they are trained from (O’Neil 2017). Second, people may consider that machines are more homogeneous than humans—that is, that a machine being biased is a sign that all machines are comparably biased (Longoni et al. 2022). Third, people may expect machines to inherit from humans not only biases but also the difficulty of fixing them, perhaps underestimating our ability to reprogram machines given our relative inability to reprogram humans (Mullainathan 2019).

Experimental results, however, suggest that people do not feel especially outraged when machines discriminate, or at least not as outraged as they would feel if humans discriminated (Bigman et al. 2023, Hidalgo et al. 2021). There is also a growing body of evidence suggesting that the very groups that feel at risk of biased human decisions may be the least averse to letting machines make decisions—seemingly because they are worried enough about the current decisions of humans to be willing to take a chance with machines (Bigman et al. 2021; Fumagalli et al. 2022; Jago & Laurin 2022; Pammer et al. 2021, 2023).

If these results are confirmed, they may create conflicts about how best to listen to the voices of the groups who are currently experiencing discrimination. When making the decision to deploy implicit moral machines, it is ethical to take into account the preferences of the persons who might be adversely and disparately impacted by the machines and to trust their lived experience of discrimination. In the context of algorithmic decisions, however, we may also need to be mindful of the knowledge that nonexperts have acquired and consider whether this knowledge is sufficient to express an informed opinion. In this space, there is a great need for clear, interactive simulations and visualizations that can help people choose their own algorithm and get first-hand experience of how implicit moral machines may affect them (Hao & Stray 2019).

**2.1.3. Blame and other reactions to harm.** So far we have focused on people’s requirements and expectations when it comes to letting implicit moral machines make consequential decisions. We now turn to people’s reactions when machines do not meet their expectations, compared to their reactions when human agents do not meet expectations. When human agents make harmful mistakes, other humans experience a manifold of negative reactions about the agent. Depending on how bad the mistake was, whether it was preventable, and whether it might have been intentional, people experience emotions such as anger and outrage, place responsibility and blame on the agent, and consider whether to punish the agent or terminate their employment (Cushman 2015, Malle et al. 2022). Do they experience the same emotions about machines, and if so, to a greater or

lesser extent? It may seem bizarre, from a rational perspective, to be angry at a machine, to hold it responsible, or to blame it for the outcome of its decision: Our anger means nothing to machines, nor our punishments. However, from a psychological perspective, people do seem to experience toward machines the same manifold of negative reactions they experience toward humans, perhaps because the machines are perceived as autonomous enough to warrant these reactions (Bigman et al. 2019, Epstein et al. 2020, Franklin et al. 2022).

In fact, when implicit moral machines make mistakes, people may experience stronger reactions than when humans make comparable mistakes. This phenomenon is clear in the domain of automated driving, across many experiments comparing people's reactions to crashes caused by human drivers and their reactions to crashes caused by self-driving cars (Franklin et al. 2021, Hidalgo et al. 2021, Hong et al. 2020, Liu & Du 2022, Liu et al. 2019a). All other things being equal, people judge crashes as more severe and less acceptable when they are caused by self-driving cars and place more blame and responsibility on a self-driving car causing a crash than on a human causing a comparable crash. It is not clear yet whether this pattern generalizes to other domains (Lima et al. 2021, Srinivasan & Sarial-Abi 2021). In particular, we already mentioned that people experience stronger negative reactions when humans discriminate than when machines do the same (Bigman et al. 2023, Hidalgo et al. 2021)—perhaps because people are angry at the idea that human discrimination may be intentional, while they do not hold the same suspicion about machine discrimination.

While it is theoretically and methodologically interesting to compare the blame incurred by humans and machines that make the same mistake, it may be more realistic to investigate situations in which human and machine jointly produce a mistake. Indeed, there may not be many situations (other than fully autonomous driving) where machines are allowed to make dangerous decisions without any human supervision. Since there will almost always be a human in the same loop as the machine, mistakes will most often be the result of a joint failure of human and machine. So, how do people allocate responsibility and blame between humans and machines, when both contributed to a harmful mistake?

Once more, the bulk of the available evidence comes from the domain of automated driving. Recall that people were less severe toward humans who caused a crash than toward machines that caused a comparable crash. Remarkably, this pattern reverses when human and machine jointly produce a crash (Awad et al. 2020b, Beckers et al. 2022, Liu et al. 2021, Wotton et al. 2022). For example, when a semiautonomous vehicle and its human driver in the loop both fail to steer away from a pedestrian, people typically blame the human the most for the resulting collision. It is not yet clear why people blame machines more than humans when they fail alone, only to blame humans more than machines when they fail together. In any case, it would be useful to collect data in other domains than self-driving cars in order to assess the transportability of this blame reversal effect (Shank et al. 2019).

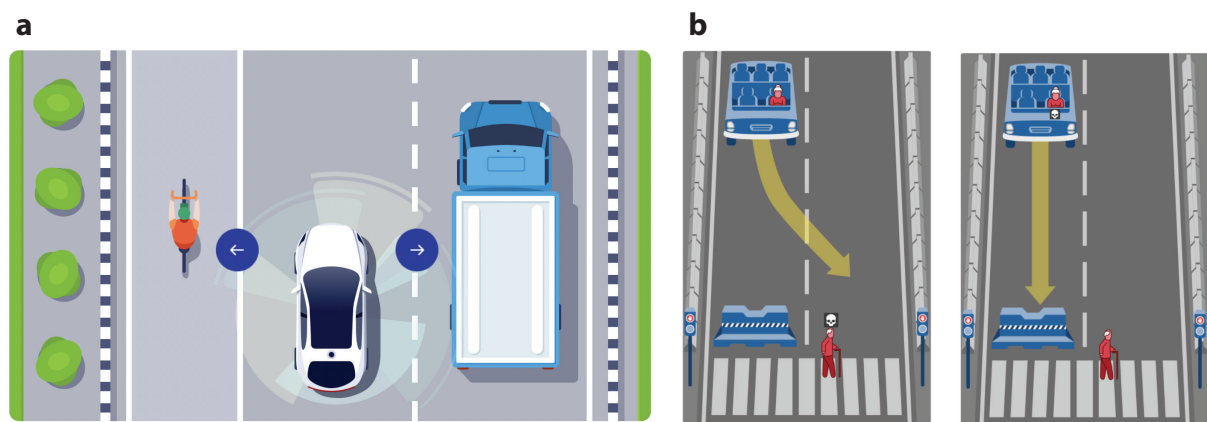
## 2.2. Explicit Moral Machines

Implicit moral machines do not attempt to solve moral dilemmas; explicit moral machines do. Indeed, explicit moral machines either solve moral dilemmas as their main function or are susceptible to encounter moral dilemmas in some situations and must accordingly be equipped to solve these dilemmas when they arise. Some moral dilemmas take the form of a conflict between two ethical principles. For example, a machine that performs content moderation online, at a scale or speed that prevents continuous human oversight, may have to routinely arbitrate between the value of free speech and the duty to suppress offensive or harmful content. In another context, a medical AI may have to arbitrate between immediately providing its best diagnosis even though it cannot explain its reasoning to humans or recommending further tests to improve explicability,

at the risk of delaying a time-sensitive diagnosis. It is very common in AI ethics to provide lists of moral values or ethical principles that AI should simultaneously pursue (e.g., beneficence, privacy, dignity, and transparency), but it is far less common to provide guidance on what machines should do when these values are in conflict (Mittelstadt 2019, Morley et al. 2020). One reason is that broad ethical principles such as dignity and privacy are hard to quantify, making it difficult to operationalize their trade-offs in policy guidelines as well as in experimental work (but see Kozyreva et al. 2023, Nussberger et al. 2022). Perhaps as a result of this difficulty, the psychological literature on explicit moral machines has mostly focused on another kind of dilemma, one that seems more amenable to experimental investigation.

This second kind of moral dilemma typically concerns the allocation of a scarce resource, with detrimental consequences for the humans who are unprioritized in the allocation decision. Consider, for example, the problem of kidney paired donation (Freedman et al. 2020). A large share of kidney transplants involve a living donor, who is usually a spouse or a relative of the candidate, but all too frequently, the potential donor is a poor match for the candidate they volunteered to help. In such a situation, one solution is to enter all candidates and prospective donors in a database, which is then fed to an algorithm that seeks two-way, three-way, or even more complex chains of donations, so that as many candidates as possible find a compatible donor. This algorithm does not simply seek to maximize the number of donations, though, but it also uses a complex priority scheme that balances many factors such as the age of the candidates, how long they have been registered in the program, their travel distance to the transplantation center, or their baseline chance to find a donor in the general population. The machine must engage in trade-offs between all these factors in order to decide who will receive a kidney and who will remain on the waiting list.

Consider now the example of autonomous vehicles (AVs), in which the scarce resource to allocate is road safety. As implicit moral agents, AVs are expected to lead to an absolute increase in road safety; but AVs are also explicit moral agents in the sense that every action they take can redistribute relative levels of safety among the road users that surround them (Bonnefon et al. 2019, Goodall 2016). This is illustrated in **Figure 1a**, in which the lateral position of the AV redistributes relative safety among the cyclist to its left, its own passengers, and the truck driver to its right. In a more extreme example (**Figure 1b**), a collision is unavoidable, and the AV must



**Figure 1**

Examples of dilemmas faced by autonomous vehicles (AVs) as explicit moral agents. (a) Depending on its lateral positioning, the AV redistributes safety among the cyclist, the AV's passenger, and the truck driver. (b) In case a collision is unavoidable, the AV may have to decide whom to save—for example, its passenger or a pedestrian.

## BLAMING MACHINES FOR SOLVING DILEMMAS

Humans will be blamed for the way they solve a moral dilemma, whatever they do. The same holds for machine, only with a twist: The blame incurred by a machine may not be distributed across possible decisions the same way it is distributed across human decisions. For example, in classic dilemmas in which an agent must decide whether to save several lives by sacrificing one, humans are blamed more when they choose to sacrifice one, but this pattern is eliminated or even reversed when a machine makes the decision (Komatsu et al. 2021, Malle et al. 2015). This implies that delegating difficult moral decisions to machines may not only remove a psychological burden from humans but also change social expectations about which decision should be made (Gill 2020).

decide whether to save its passenger or a pedestrian (Bonnefon et al. 2016). In both situations, the AV must be endowed with the ability to make a moral calculation about whose safety should take priority.

While people are often uncomfortable with the idea of letting machines make moral decisions (Bigman & Gray 2018, Dietvorst & Bartels 2022, Shariff et al. 2017), there is a case to be made that it is good to let machines solve moral dilemmas, even and especially when their decisions have unavoidably tragic consequences. If we agree that making such decisions inflicts an emotional cost on the decision maker, in both the short and the long term, we may agree that it is good to delegate this burden to machines, which do not experience psychological suffering (Danaher 2022). In like vein, we know that humans will inevitably be blamed for the way they solved a moral dilemma, since by definition a moral dilemma has no universally accepted solution—hence, we may want to relieve human decision makers from unavoidable blame by delegating the decision to a machine (see the sidebar titled *Blaming Machines for Solving Dilemmas* for further discussion).

If explicit moral machines are to make moral trade-offs, we need to provide them with the goals and priorities they should pursue. This challenge is part of the value alignment problem (Gabriel 2020): To ensure that machines solve moral dilemmas in a way that is compatible with the goals and priorities of humans, we need to know what these human goals and priorities are and to find a way to teach them to machines. Here we are concerned with the first objective, which falls squarely within the purview of moral psychology. We consider in turn some specific difficulties that moral psychologists face when collecting human moral preferences for the purpose of teaching them to machines: who to ask, how to ask, and what to do with the answers.

**2.2.1. Value alignment: who to ask.** Not everyone agrees about what should be done in a moral dilemma, or what values should take priority in a moral trade-off. So, who do we ask for their moral preferences, when we want to inform the decisions of machines? A good place to start is to ask ethicists, who are trained to think about these issues and have a deep understanding of their implications. Ethicists, however, are not immune to biases (Schwitzgebel & Cushman 2015), and they do not always come to an agreement—for example, a German national ethics committee could not reach a consensus on what an AV should do when deciding whether to save its passengers or other road users (Luetge 2017). There is another expert group we can ask for their preferences, namely, the people who build the machines and who have a detailed understanding of what AI can actually do. For example, we could ask the AV industry what they believe AVs should do when deciding to save their passengers or other road users. One problem is that the industry is very reluctant to engage in this debate (Martinho et al. 2021), since any position they take may alienate either their consumer base or the general population. Experts from the AV industry may also feel a duty to protect their customers, which could explain why they have a stronger preference to save passengers compared to the general population (Zhu et al. 2022).



Asking AI developers and ethicists for their informed preferences is important in view of their expertise, but we also need to document the preferences of the laypersons who will adopt the technology. Consider again the dilemma of an AV that needs to decide whether to prioritize the life of its passengers or that of other road users. However rare this dilemma might be, it weighs heavily on the minds of consumers, to the point of being cited as one of the top issues that will determine their decision to adopt AVs (Gill 2021). In this context, learning about consumers' preferences is not merely a marketing exercise. The main promise of AVs is that they can reduce the number of road casualties by being safer than human drivers; but these lives will not be saved if consumers opt out of the technology because they are unsatisfied with or even outraged by the way AVs solve moral dilemmas (Bonneton et al. 2020b, De Freitas & Cikara 2021). As a result, learning the moral preferences of consumers may be a prerequisite for explicit moral machines to deliver their benefits. Explicit moral machines do not impact only the outcomes of their adopters, though. By design, they can create externalities for other stakeholders. For example, AVs do not merely affect the safety of their passengers but also distribute risk to all road users around them. As a result, other road users (including pedestrians) should be given a voice when collecting preferences about the moral priorities that determine the behavior of AVs.

In sum, value alignment requires the collection of human moral preferences to inform the behavior of explicit moral machines—a process that requires a decision about whose values will be collected and how they should be weighted when different groups have different preferences. These normative issues are complex, but they arguably fall beyond the purview of moral psychology. Moral psychologists have an important role to play, however, in bringing their expertise to the matter of how best to measure moral preferences about explicit moral machines.

**2.2.2. Value alignment: how to ask.** Measuring moral preferences is never easy, and measuring preferences about explicit moral machines comes with its own set of challenges. First, explicit moral machines may need to balance a great number of conflicting values or priorities. For example, kidney paired donation algorithms may balance up to a dozen priorities, including the quality of the match between donor and candidate, the statistical rarity of potential donors for a given candidate, the age of the candidate at registration in the program as well as their waiting time in the program, the blood types of donor and candidate, and the possibility that the candidate has donated a kidney themselves. When distributing risk around them, AVs may need to consider the number of potential victims, their mode of transportation, their age, whether they are currently on the road or the sidewalk, and yet other variables. The high-dimensional nature of these choices may lead to an exploding number of experimental treatments, resulting in a need for an unpractical number of research participants. The Moral Machine experiment (Awad et al. 2018) considered nine possible priorities for AVs to decide which group of road users to save or to sacrifice, which led to millions of possible scenarios. Exploring this enormous space was only possible because the experiment went viral, collecting data from millions of participants. Not every experiment can go viral, though, which means that moral psychologists have difficult choices to make when deciding how complex they want their scenario space to be.

Many other design choices will impact the feasibility of such experiments and the interpretation of their results. For example, the Moral Machine experiment purposefully used stylized scenarios when a collision is unavoidable (Awad et al. 2020a), but more realistic scenarios would have manipulated the probability of the collisions (Krügel & Uhl 2022). Participants were asked what the AV should do, but other questions can lead to different results—for example, asking participants what AV behavior they would prefer as passengers (Bonneton et al. 2016, Liu & Liu 2021, Takaguchi et al. 2022), as other road users (Martin et al. 2021, Mayer et al. 2021), or from under a veil of ignorance (Huang et al. 2019). Other experiments may opt out of asking

participants to state their preferences and try instead to reveal their preferences by placing them in a virtual environment where they need to make themselves the same moral decisions that AVs will face (Faulhaber et al. 2019, Samuel et al. 2020). Given the relative novelty of explicit moral machines as a topic of investigation for moral psychology, the field may be best served by embracing this diversity of methods and designs in order to build a comprehensive description of the moral values that people may want to see embedded in machines. This inclusive approach is especially important in view of what we will do with these data, as we discuss in Sections 3 and 4.

**2.2.3. Value alignment: what to do with the answers.** It seems consensual to say that moral psychology has an important descriptive role in documenting the values and priorities that laypersons would want explicit moral machines to pursue (Awad et al. 2022). What is much more controversial is to decide what prescriptive weight these data should have in the policies that will regulate the behavior of the machines. Clearly, no one wants these policies to be driven solely by the preferences of laypersons—but should these preferences be discarded entirely?

A promising approach to this question is to jointly consider the degree of consensus or division among experts and the degree of consensus or division among laypersons (Savulescu et al. 2021). Consider first the situation where experts show strong consensus about what a machine should do. If laypersons show the same consensus, the case is closed. If laypersons are divided about what the machine should do, then the proper course of action is probably to follow the expert consensus while building up a strong and clear case for this consensus in terms that the public can understand. If laypersons show a strong consensus against the consensus of the experts, the situation is more difficult, but it is also possible that the public consensus is based on bias more than reason, which is something that moral psychologists are equipped to show.

Consider now the situation where experts themselves are divided, and this division reflects a reasonable moral disagreement. In this case, it may be appropriate to follow the public consensus, if there is one. However, this requires researchers to be very careful about establishing this consensus and to make sure it does not reflect, for example, the biases and prejudices of the majority. This is why we believe it is especially important for moral psychologists to explore an exhaustive range of methods and controls in order to make sure that the public consensus is robust across experimental designs and demographics, as well as free of prejudice and bias, before it is allowed to arbitrate over the disagreements of experts.

### 3. MACHINES AS MORAL PATIENTS

So far, we have considered situations where machines are (implicit or explicit) moral agents, that is, situations in which machines perform actions whose consequences affect people. We now flip the table and consider situations in which machines are moral patients, that is, situations in which people perform actions that affect machines. This may sound strange, since machines have no affects nor needs or desires for anything. Even though people are well aware of this, they can still feel empathy for machines (see the sidebar titled *Empathy for the Machine*) or consider that machines want things in a certain sense, things that can be given or denied. In other terms, and as we consider in more detail in the rest of this section, people sometimes assume that machines have preferences, and this can turn machines into moral patients that experience preferred or dis-preferred outcomes as a result of the actions taken by other agents (Pauketat & Anthis 2022).

This is especially important when people have an opportunity to cooperate with a machine. Cooperative interactions with intelligent, autonomous machines are not yet a common experience, but this is likely to change in the future. Cooperation with machines is already a reality in industry settings (Villani et al. 2018), and soon enough, road users will have to cooperate with AVs to make

## EMPATHY FOR THE MACHINE

While people understand that robots do not experience physical pain or psychological distress, they can nevertheless feel emotionally uncomfortable when humans direct toward robots the kind of behavior that would qualify as abuse if directed toward other humans. For example, research participants show physiological signs of discomfort when watching a baby dinosaur robot being punched and choked (Rosenthal-von der Pütten et al. 2013) or a robot hand being cut by a knife (Suzuki et al. 2015); they hesitate when asked to strike a robot (Darling et al. 2015) or to topple a block tower that a robot built and pretends to care about (Briggs & Scheutz 2014); and they are likely to ask a research confederate to stop when they see the confederate insulting and roughing up a robot (Connolly et al. 2020).

traffic safe for everyone (Schwartz et al. 2019). Many participants in online communities or social networks already have with bots the same kind of cooperative (or uncooperative) interactions that they have with humans (Seering et al. 2018, Shao et al. 2018, Stella et al. 2018, Tsvetkova et al. 2017): People and bots can retweet or block one another; Reddit users sometimes congratulate bots for good behavior, but sometimes they report them to moderators; and Wikipedia editors can cooperate with bots on an article or engage in an editing war against them. As interactions with intelligent machines become more commonplace, how will humans and machines initiate and sustain cooperation?

Mutually beneficial cooperation between humans often relies on a positive concern for the outcomes of others—a preference for the satisfaction of the preferences of others. Cooperation is easier if my other-regarding preferences are prosocial, that is, if I derive some measure of satisfaction from doing good to others. Conversely, cooperation is usually more difficult if my other-regarding preferences are antisocial or even just callous—that is, if I derive satisfaction from doing ill to others, or if I am entirely indifferent about what happens to others and only care about my own outcomes. What happens when humans have an opportunity to cooperate with machines? What are their machine-regarding preferences? This is the topic of the next section.

### 3.1. Machine-Regarding Preferences

Cooperation between humans does not necessarily involve money. People can volunteer their time and skills to help others, provide advice, share tools, advocate for a cause, or donate blood. While it is possible to study all these currencies in behavioral experiments that investigate cooperation, experimental economics has popularized the assumption that it is possible to capture the manifold of human cooperation by using lab-based games with financial incentives, such as dictator games, prisoner's dilemmas, ultimatum games, or public good games. Incentivized games provide a controlled, stylized environment to measure other-regarding preferences and prosocial behavior that allows for easy comparison of studies and experimental treatments. As a result, many studies of human-machine cooperation have used the same games, only replacing some human players by intelligent machines, in order to document changes in human behavior when playing incentivized games with machines rather than humans (March 2021). These studies carry over the assumption that just as money can be used as a proxy for the many currencies of human-human cooperation, so it can be used as a proxy for the many currencies of human-machine cooperation. In the rest of this section, we proceed with this assumption—but see the sidebar titled *What Do Machines Do with Money?* for a closer examination.

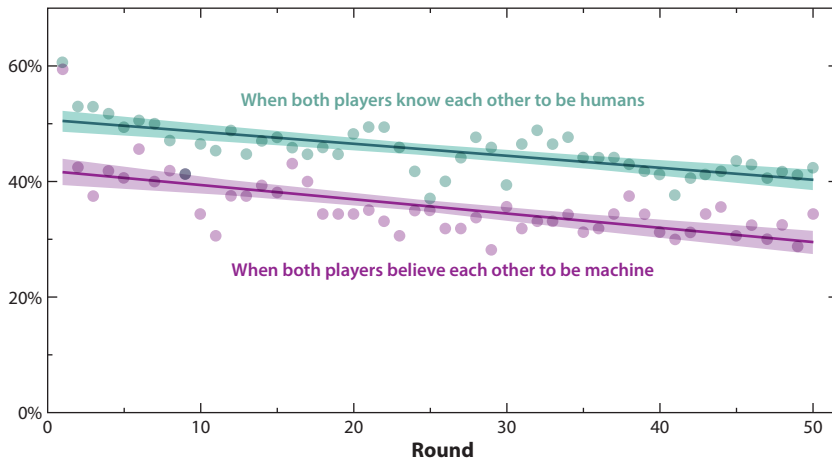
Findings on human-machine cooperation in incentivized games show remarkable convergence. In a nutshell, people do show some measure of prosocial machine-regarding preferences,

## WHAT DO MACHINES DO WITH MONEY?

If you had to split some money between yourself and, say, a tree, you would probably wonder about what happens to the money you give to the tree, since trees have no use for money. The same question holds in experiments in which people share money with machines or help machines make money. Presumably, what truly happens in most cases is that the money earned by the machine goes back to the research fund of the experimenters—but this is not usually made clear to research participants. Indeed, in a survey of 160 experiments, von Schenk et al. (2022) observed that 82% of instructions did not give any explanation about what happened to machine earnings. (The rest was split between pretending that machines would keep the money, reminding that machines had no use for money, and explaining that the machine earnings would be transferred to a human.) So, if people wonder what machines could use money for, and if experimenters have no answer to offer, is money the right currency to study human–machine cooperation? There are two arguments for believing so. First, it is not like any other currency would be better, since machines do not care about anything in the sense of feeling a desire or a need for something. Second, people seem to agree that machines still want money, in the sense of being programmed to do so, to the same extent that they want retweets or other cooperative currencies used in online communities (Makovi et al. 2023). As a result, money in incentivized games seems an acceptable proxy for the currencies used in real-life human–machine cooperation.

and cooperation does not disappear when humans play with machines—but it does not reach the level of human–human cooperation. In other words, all findings suggest the existence of a machine penalty in cooperative games. For example, in a one-shot trust game, human second-movers expected the same level of cooperation from human and machine first-movers, but only 34% reciprocated the trust of a machine, compared to 75% who reciprocated the trust of a human; likewise, in a one-shot prisoner’s dilemma, people expected the same level of cooperation from humans and machines but cooperated with only 36% of machines compared to 49% of humans (Karpus et al. 2021). In a one-shot dictator game, people allocated 39% of their endowment to a human but only 16% to a machine; and in a one-shot public good game, people contributed about 55% of their endowment to the common pool when playing with humans but only 40% when playing with machines (Nielsen et al. 2022b).

One-shot games thus suggest that people do not initiate cooperation with machines to the same level they initiate cooperation with humans: Cooperation does not drop to zero, but it suffers from a machine penalty. Repeated games allow to study the dynamics of the machine penalty and its evolution through repeated interaction (Crandall et al. 2018). Findings suggest that the machine penalty carries unchanged over repeated interactions, but their interpretation can be complicated by the fact that in repeated games, human decisions can be impacted by the strategy chosen by the machine, which may be different from the strategies commonly adopted by humans (Sandoval et al. 2016). One way to address this difficulty is to use deception, that is, to pair players with humans they believe to be machines, or to pair them with machines they believe to be humans. Such deception allows to measure the mere effect of believing that one’s partner is human or machine, independently of the strategy adopted by the partner. **Figure 2** displays the results of one such experiment (Ishowo-Oloko et al. 2019), in which human players either knew each other to be humans or believed each other to be machines. As is common with repeated prisoner’s dilemmas, cooperation steadily decreases over time when both players know each other to be humans. When both players believe each other to be machines, the negative dynamics is very similar, and the machine penalty carries over unchanged over time.



**Figure 2**

Over 50 rounds, cooperation between two human players in a prisoner’s dilemma steadily decreases. The dynamics are the same when the players believe each other to be machines, hence the parallel regression lines, but the machine penalty carries over time, hence the vertical distance between the two lines. Data replotted from the source file of Ishowo-Oloko et al. (2019).

### 3.2. Overcoming the Machine Penalty

The machine penalty is not only a phenomenon we need to understand but also, arguably, a problem we must solve. For the last 20 years, celebrated milestones in AI research were often tied to competition against humans—be it when IBM Watson defeated the two highest-ranked *Jeopardy!* players (Ferrucci et al. 2010) or when DeepMind’s AlphaGo defeated the top Go player Lee Sedol (Silver et al. 2016). While surpassing human performance is an important goal of AI (in particular when it behaves as an implicit moral agent), there is an increasing recognition that in order to fulfill the true potential of AI, we need to put as much effort in human–AI cooperation as in human–AI competition. There is a technical side to this challenge, since it may require designing AI systems that understand and respond appropriately to human intentions and goals (Dafoe et al. 2021); but there is also a psychological side to the challenge, which requires understanding why humans are reluctant to cooperate with machines and designing interventions that can overcome this machine penalty.

It is perhaps natural to start with interventions that give machines more human-like traits. After all, if people do not cooperate with machines as much as they do with humans, perhaps we can narrow the gap between cooperation rates by making machines look or feel more like humans. This humanization strategy may help people activate with machines the same cooperation templates they activate with humans or the frames of reference they use to interpret the behavior of cooperation partners, thus increasing their trust and comfort in this new situation (Nielsen et al. 2022a). We may then expect the humanization strategy to increase in efficacy when machines are humanized to a large degree compared to when machines are humanized to a minimal degree. Experimental findings, however, tell a more complicated story.

Minimally humanized robots typically fail to elicit more cooperation than nonhumanized robots. Examples of minimal humanization include giving the robot an ovoid shape augmented with eyes as compared to an insectoid appearance (De Kleijn et al. 2019), or endowing a non-humanoid robot with some emotional displays, such as stylized angry, sad or happy eyes, and recorded sighs and laughter (Hsieh & Cross 2022). These experiments do not report significant

effects on cooperation, suggesting that minimal humanization is insufficient to overcome the machine penalty. Climbing up the humanization gradient does not improve cooperation much, and it can even make things worse due to the uncanny valley effect—that is, the feeling of strangeness and discomfort elicited by a machine that is largely but not quite human-like. For example, a study using 80 robotic faces going from entirely machine-like to entirely human-like found that cooperation was at its lowest for machines that were placed at two-thirds of the humanization gradient (Mathur & Reichling 2016), and other studies showed that even more human-like robots failed to eliminate the machine penalty (Cominelli et al. 2021, Złotowski et al. 2016). Intriguingly, the few studies that succeeded in reducing the machine penalty through (moderate) humanization did so by gendering the machine as female, through stylized cues such as suggestions of long hair or breasts (Bernotat et al. 2021, Eyssel & Hegel 2012). While this strategy may indeed prove somewhat efficient at reducing the machine penalty, it seems ethically problematic to exploit and perpetuate gender stereotypes about women being less competitive or more nurturing, just as it would seem problematic to systematically give AI assistants a female voice (Fossa & Sucameli 2022).

All humanization strategies we reviewed so far were nondeceptive, in the sense that while the machine was made more human-like, it was never described as being human. If we remove that constraint, we reach the highest possible level of humanization, machines that pretend to be humans. This is done easily enough in most experimental protocols that use incentivized games, since these protocols are usually designed to remove all visual or verbal interactions between players. If players are identified with headshots, machines can create synthetic faces for themselves, which can be both realistic and especially trust inducing (Nightingale & Farid 2022). Unsurprisingly, this deceitful form of humanization eliminates the machine penalty (Ishowo-Oloko et al. 2019): If people do not know they are cooperating with machines, they do not manifest the machine penalty. Once more, though, this solution creates ethical issues, since AI codes of ethics typically emphasize that machines should never be allowed to pass as humans (O’Leary 2019).

In sum, humanization strategies usually fail to reduce the machine penalty, and the ones that succeed (partially or totally) fall short of current ethical standards. As a result, there is a need for further research that would seek to improve human–machine cooperation without resorting to the humanization of machines. One promising direction may be to embrace the fact that intelligent machines are newcomers in our social and cooperative interactions and to accept that dealing with these newcomers may require new social norms (Makovi et al. 2023). In other words, rather than making machines more human-like in the hope that people will apply to them the old social norms they apply to humans, we could experiment on the new social norms that will develop around the new entrants in our social world, intelligent machines.

## 4. MACHINES AS MORAL PROXIES

By design, AI enables machines to make autonomous decisions on behalf of human stakeholders. This raises the possibility of delegating unethical behavior in a way that distances the human from the act. AI also offers a further possibility, namely that of mediating human communication in a morally relevant manner. We explore each of these possibilities in turn.

### 4.1. Delegation to Machines

People delegate a growing number of tasks to AI agents (de Melo et al. 2018). Current and near-term possibilities are as diverse as setting prices in online markets (Calvano et al. 2020a), interrogating suspects (McAllister 2016), and marketing to customers (Cheng & Jiang 2022). This creates many opportunities to delegate unethical behavior to machines.

First, AI can be used by people who have malicious intentions to scale up criminal or unethical behavior. Recent advancements in deep learning, and specifically generative adversarial networks (GANs), have made it easier to create fake content that looks genuine (Caldwell et al. 2020). Those who have malicious intentions can benefit from using AI hench-agents because AI can act independently and has the potential to cause harm with unparalleled efficiency and at scale. Moreover, these AI hench-agents may be harder to trace back to the original source. AI-powered deepfakes can create fake identities, which allows phishing attacks to become more personalized and effective. These attacks, also known as spear phishing (Seymour & Tully 2016), put a new spin on identity theft and can have devastating results (Jagatic et al. 2007). Reflecting on this emerging worry, a panel of experts has nominated deepfakes as the most dangerous tool for AI-enabled crime (Caldwell et al. 2020).

Delegation of criminal or ethically questionable behavior to AI agents might be attractive for reasons other than scalability. When people delegate tasks to AI agents, this creates a combination of psychological factors that can lead to unethical behavior, such as anonymity (Ostermaier & Uhl 2017), psychological distance from victims (Köbis et al. 2019), and undetectability (Hancock & Guillory 2015, Rauhut 2013). The often-incomprehensible workings of algorithms create ambiguity (Miller 2019). Letting such black box algorithms execute tasks on one's behalf increases plausible deniability and obfuscates the attribution of responsibility for the harm caused. If any harm does become apparent, blame and responsibility can be deflected to the delegate, which may alleviate the (legal or psychological) guilt experienced by the remitter. Indeed, people tend to prefer delegation, even if it entails explicit instruction to break ethical rules, such as when using hench-persons (Drugov et al. 2014).

Ambiguity is another mechanism through which unethical behavior can be delegated to machines. More often than not, people do not explicitly instruct their delegates to break ethical rules but instead merely define their desired outcome and turn a blind eye to how it is achieved. By doing so, the remitter avoids direct contact with the victim and can willfully ignore, through deliberate ignorance (Hertwig & Engel 2016), any possible ethical rule violations that may occur as a result of the delegation (Drugov et al. 2014, Van Zant & Kray 2014).

Delegation to AI may also cause moral violations without any bad intent (Thomas et al. 2019). For example, someone may use algorithmic prices to sell goods on online markets, without being aware that algorithms might coordinate and set collusive prices (Calvano et al. 2020b, Wellman & Rajan 2017). Marketers who rely on AI-powered sales strategies might be unaware of the fact that the AI agent employs deceptive tactics to reach sales goals.

Not all delegation is bad, of course. One may indeed delegate morally desirable actions to AI agents. Specifically, delegating morally desirable actions such as charitable donations to an AI agent may act as a commitment device (Bryan et al. 2010) that increases the magnitude and frequency of such actions. There are also opportunities to delegate an advisory role to AI agents, enabling them to dynamically suggest moral behavior to the human (Giubilini & Savulescu 2018).

## 4.2. Machine Masquerade

We close this article with our shortest and most speculative section. So far we considered the possibility for people to send a machine proxy in a moral interaction, in the sense that they delegate their decisions to the machine. In this final section, we consider the possibility for people to participate themselves in an interaction, only under a disguise provided by the machine. Under this machine masquerade (known as AI-mediated communication; Hancock et al. 2020), people use technology to modify the way they write, talk, or look in order to change the behavior of their partners. Moral psychology has given little attention so far to AI-mediated communication, but

this is likely to change given the incoming availability of machine masquerade tools, the way they will transform moral interactions, and the ethical challenges they raise.

Many people are already familiar with machine-generated replies to text messages or emails as well as image filters that improve the appearance of the subject; but AI is poised to allow them much more powerful and flexible forms of transformation. Written text, profile pictures, and voice and facial dynamics in live online interactions can already be altered to achieve various presentation goals. While not everyone will have immediate access to all these technologies (Goldenthal et al. 2021), their adoption can be very fast. Consider the case of OpenAI's ChatGPT, as of the writing of this article a state-of-the-art language model with a user-friendly interface that allows people to easily experiment with various prompts and requests. Within weeks of its public launch, ChatGPT attracted more than a hundred million users, affording them seemingly endless possibilities. Students could use ChatGPT to sound more competent, business owners could use it to sound more trustworthy, and social media users could ask it to generate posts in line with the image they wished to project or the moral virtues they wished to signal.

We know very little yet about how people will seize and judge these opportunities, at which scale, and to which effects. People who use machines to write for them are perceived as less trustworthy, according to studies using hypothetical emails (Liu et al. 2022), hypothetical Airbnb profiles (Jakesch et al. 2019), and actual text conversations (Hohenstein et al. 2023); but there is much more to be done to understand the material and reputational benefits that people can achieve if their use of machines is not discovered, and how much of these benefits are conserved depending on the way the use of machines is disclosed. Compare, for example, a social media user who is posting content that they secretly asked a machine to generate in order to signal a commitment to gender equality and a social media user who is disclosing on their profile that they are systematically asking a machine to alter their posts in order to remove gender biases. Such scenarios are no longer far-fetched, and we need moral psychology to understand the effects they will have as well as the reactions they will trigger.

Machine masquerade is not restricted to written text; it can also alter the way we look and the way we sound. People can already experiment with generative AI to create their profile pictures, and the technology to alter voices is already perfected. This means that people can ask machines to alter their face in order to appear more dominant or more trustworthy, or to alter their voice to sound more articulate or more cheerful (Gueraouaou et al. 2022). These alterations can change the outcomes of moral interactions, but they can also raise new ethical issues for moral psychology to investigate, such as the conflict between reducing discrimination and jeopardizing inclusion. For example, machines can remove the foreign accent of call center employees, which decreases the likelihood they will receive racist abuse from angry customers—but this can be construed as a step in the wrong direction, as it would amount to considering that the fix to racism is not to reduce prejudice but to accommodate it by whitening the voice of its victims (McCallum & Vallance 2022).

In sum, machine masquerade offers a vast new field of investigation for moral psychology, aimed at understanding how people will use technology to alter their presentation, either for the purpose of changing the outcomes of moral interactions or for the purpose of managing their moral reputation; how this processes may be moderated by different forms of disclosure; and how society will deal with the new ethical dilemmas raised by this technology.

## 5. CONCLUSION

We have not addressed every issue at the intersection of AI and moral psychology. Questions about how people perceive AI plagiarism, about how the presence of AI agents can reduce or enhance



trust between groups of humans, and about how sexbots will alter intimate human relations are the subjects of active research programs. Many more yet unasked questions will only be provoked as new AI abilities develop. Given the pace of this change, any review paper will only be a snapshot. Nevertheless, the very recent and rapid emergence of AI-driven technology is colliding with moral intuitions forged by culture and evolution over the span of millennia. Grounding an imaginative speculation about the possibilities of AI with a thorough understanding of the structure of human moral psychology will help prepare for a world shared with, and complicated by, machines.

## DISCLOSURE STATEMENT

The authors are not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

## ACKNOWLEDGMENTS

J.-E.B. acknowledges support from grant ANR-19-PI3A-0004, grant ANR-17-EURE-0010, and the research foundation TSE-Partnership. A.S. acknowledges support from a Canada 150 Research Chair grant from the Social Sciences and Humanities Research Council of Canada.

## LITERATURE CITED

- Angwin J, Larson J, Mattu S, Kirchner L. 2016. Machine bias. *ProPublica*, May 23. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- Awad E, Dsouza S, Bonnefon JF, Shariff A, Rahwan I. 2020a. Crowdsourcing moral machines. *Commun. ACM* 63(3):48–55
- Awad E, Dsouza S, Kim R, Schulz J, Henrich J, et al. 2018. The Moral Machine experiment. *Nature* 563(7729):59–64
- Awad E, Levine S, Anderson M, Anderson SL, Conitzer V, et al. 2022. Computational ethics. *Trends Cogn. Sci.* 26(5):388–405
- Awad E, Levine S, Kleiman-Weiner M, Dsouza S, Tenenbaum JB, et al. 2020b. Drivers are blamed more than their automated cars when both make mistakes. *Nat. Hum. Behav.* 4(2):134–43
- Beckers N, Siebert LC, Bruijnes M, Jonker C, Abbink D. 2022. Drivers of partially automated vehicles are blamed for crashes that they cannot reasonably avoid. *Sci. Rep.* 12:16193
- Bernotat J, Eyssel F, Sachse J. 2021. The (fe)male robot: how robot body shape impacts first impressions and trust towards robots. *Int. J. Soc. Robot.* 13(3):477–89
- Bigman YE, Gray K. 2018. People are averse to machines making moral decisions. *Cognition* 181:21–34
- Bigman YE, Waytz A, Alterovitz R, Gray K. 2019. Holding robots responsible: the elements of machine morality. *Trends Cogn. Sci.* 23(5):365–68
- Bigman YE, Wilson D, Arnestad MN, Waytz A, Gray K. 2023. Algorithmic discrimination causes less moral outrage than human discrimination. *J. Exp. Psychol. Gen.* 152(1):4–27
- Bigman YE, Yam KC, Marciano D, Reynolds SJ, Gray K. 2021. Threat of racial and economic inequality increases preference for algorithm decision-making. *Comput. Hum. Behav.* 122:106859
- Birhane A. 2022. The unseen Black faces of AI algorithms. *Nature* 610:451–52
- Bonnefon JF, Černý D, Danaher J, Devillier N, Johansson V, et al. 2020a. *Ethics of Connected and Automated Vehicles: Recommendations on Road Safety, Privacy, Fairness, Explainability and Responsibility*. Brussels, Belg.: Publ. Off. Eur. Union
- Bonnefon JF, Shariff A, Rahwan I. 2016. The social dilemma of autonomous vehicles. *Science* 352(6293):1573–76
- Bonnefon JF, Shariff A, Rahwan I. 2019. The trolley, the bull bar, and why engineers should care about the ethics of autonomous cars [point of view]. *Proc. IEEE* 107(3):502–4
- Bonnefon JF, Shariff A, Rahwan I. 2020b. The moral psychology of AI and the ethical opt-out problem. In *The Ethics of Artificial Intelligence*, ed. SM Liao, pp. 109–26. Oxford, UK: Oxford Univ. Press
- Bono T, Croxson K, Giles A. 2021. Algorithmic fairness in credit scoring. *Oxf. Rev. Econ. Policy* 37(3):585–617

- Briggs G, Scheutz M. 2014. How robots can affect human behavior: investigating the effects of robotic displays of protest and distress. *Int. J. Soc. Robot.* 6(3):343–55
- Bryan G, Karlan D, Nelson S. 2010. Commitment devices. *Annu. Rev. Econ.* 2:671–98
- Buolamwini J, Gebru T. 2018. Gender shades: intersectional accuracy disparities in commercial gender classification. *PMLR* 81:77–91
- Cadario R, Longoni C, Morewedge CK. 2021. Understanding, explaining, and utilizing medical artificial intelligence. *Nat. Hum. Behav.* 5(12):1636–42
- Caldwell M, Andrews JT, Tanay T, Griffin LD. 2020. AI-enabled future crime. *Crime Sci.* 9:14
- Calvano E, Calzolari G, Denicolò V, Harrington JE Jr, Pastorello S. 2020a. Protecting consumers from collusive prices due to AI. *Science* 370(6520):1040–42
- Calvano E, Calzolari G, Denicolò V, Pastorello S. 2020b. Artificial intelligence, algorithmic pricing, and collusion. *Am. Econ. Rev.* 110(10):3267–97
- Cheng Y, Jiang H. 2022. Customer–brand relationship in the era of artificial intelligence: understanding the role of chatbot marketing efforts. *J. Product Brand Manag.* 31(2):252–64
- Chouldechova A. 2017. Fair prediction with disparate impact: a study of bias in recidivism prediction instruments. *Big Data* 5(2):153–63
- Combs TS, Sandt LS, Clamann MP, McDonald NC. 2019. Automated vehicles and pedestrian safety: exploring the promise and limits of pedestrian detection. *Am. J. Prev. Med.* 56(1):1–7
- Cominelli L, Feri F, Garofalo R, Giannetti C, Meléndez-Jiménez MA, et al. 2021. Promises and trust in human–robot interaction. *Sci. Rep.* 11:9687
- Connolly J, Mocz V, Salomons N, Valdez J, Tsoi N, et al. 2020. Prompting prosocial human interventions in response to robot mistreatment. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, pp. 211–20. New York: ACM
- Crandall JW, Oudah M, Tennom, Ishowo-Oloko F, Abdallah S, et al. 2018. Cooperating with machines. *Nat. Commun.* 9:233
- Cushman F. 2015. Deconstructing intent to reconstruct morality. *Curr. Opin. Psychol.* 6:97–103
- Dafoe A, Bachrach Y, Hadfield G, Horvitz E, Larson K, Graepel T. 2021. Cooperative AI: Machines must learn to find common ground. *Nature* 593:33–36
- Danaher J. 2022. Tragic choices and the virtue of techno-responsibility gaps. *Philos. Technol.* 35:26
- Darling K, Nandy P, Breazeal C. 2015. Empathic concern and the effect of stories in human-robot interaction. In *2015 24th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, pp. 770–75. Piscataway, NJ: IEEE
- De Freitas J, Cikara M. 2021. Deliberately prejudiced self-driving vehicles elicit the most outrage. *Cognition* 208:104555
- De Kleijn R, van Es L, Kachergis G, Hommel B. 2019. Anthropomorphization of artificial agents leads to fair and strategic, but not altruistic behavior. *Int. J. Hum.-Comput. Stud.* 122:168–73
- de Melo CM, Marsella S, Gratch J. 2018. Social decisions and fairness change when people’s interests are represented by autonomous agents. *Auton. Agents Multi-Agent Syst.* 32:163–87
- Dietvorst BJ, Bartels DM. 2022. Consumers object to algorithms making morally relevant tradeoffs because of algorithms’ consequentialist decision strategies. *J. Consum. Psychol.* 32(3):406–24
- Drugov M, Hamman J, Serra D. 2014. Intermediaries in corruption: an experiment. *Exp. Econ.* 17:78–99
- Epstein Z, Levine S, Rand DG, Rahwan I. 2020. Who gets credit for AI-generated art? *iScience* 23(9):101515
- Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, et al. 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542(7639):115–18
- Eyssel F, Hegel F. 2012. (S)he’s got the look: gender stereotyping of robots. *J. Appl. Soc. Psychol.* 42(9):2213–30
- Faulhaber AK, Dittmer A, Blind F, Wächter MA, Timm S, et al. 2019. Human decisions in moral dilemmas are largely described by utilitarianism: Virtual car driving study provides guidelines for autonomous driving vehicles. *Sci. Eng. Ethics* 25:399–418
- Ferrucci D, Brown E, Chu-Carroll J, Fan J, Gondek D, et al. 2010. Building Watson: an overview of the DeepQA project. *AI Mag.* 31(3):59–79
- Fossa F, Sucameli I. 2022. Gender bias and conversational agents: an ethical perspective on social robotics. *Sci. Eng. Ethics* 28:23

- Franklin M, Ashton H, Awad E, Lagnado D. 2022. Causal framework of artificial autonomous agent responsibility. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 276–84. New York: ACM
- Franklin M, Awad E, Lagnado D. 2021. Blaming automated vehicles in difficult situations. *iScience* 24(4):102252
- Freedman R, Borg JS, Sinnott-Armstrong W, Dickerson JP, Conitzer V. 2020. Adapting a kidney exchange algorithm to align with human values. *Artif. Intell.* 283:103261
- Fumagalli E, Rezaei S, Salomons A. 2022. OK computer: worker perceptions of algorithmic recruitment. *Res. Policy* 51(2):104420
- Gabriel I. 2020. Artificial intelligence, values, and alignment. *Minds Mach.* 30(3):411–37
- Gill T. 2020. Blame it on the self-driving car: how autonomous vehicles can alter consumer morality. *J. Consum. Res.* 47(2):272–91
- Gill T. 2021. Ethical dilemmas are really important to potential adopters of autonomous vehicles. *Ethics Inform. Technol.* 23(4):657–73
- Giubilini A, Savulescu J. 2018. The artificial moral advisor: the “ideal observer” meets artificial intelligence. *Philos. Technol.* 31:169–88
- Goldenthal E, Park J, Liu SX, Mieczkowski H, Hancock JT. 2021. Not all AI are equal: exploring the accessibility of AI-mediated communication technology. *Comput. Hum. Behav.* 125:106975
- Goodall NJ. 2016. Away from trolley problems and toward risk management. *Appl. Artif. Intell.* 30(8):810–21
- Guerouaou N, Vaiva G, Aucouturier JJ. 2022. The shallow of your smile: the ethics of expressive vocal deep-fakes. *Philos. Trans. R. Soc. B* 377(1841):20210083
- Hamilton M. 2019. The biased algorithm: evidence of disparate impact on Hispanics. *Am. Crim. Law Rev.* 56:1553–77
- Hancock JT, Guillory J. 2015. Deception with technology. In *The Handbook of the Psychology of Communication Technology*, ed. SS Sundar, pp. 270–89. New York: Wiley
- Hancock JT, Naaman M, Levy K. 2020. AI-mediated communication: definition, research agenda, and ethical considerations. *J. Comput.-Mediat. Commun.* 25(1):89–100
- Hao K, Stray J. 2019. Can you make AI fairer than a judge? Play our courtroom algorithm game. *MIT Technology Review*, Oct. 17. <https://www.technologyreview.com/2019/10/17/75285/ai-fairer-than-judge-criminal-risk-assessment-algorithm/>
- Harrison G, Hanson J, Jacinto C, Ramirez J, Ur B. 2020. An empirical study on the perceived fairness of realistic, imperfect machine learning models. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 392–402. New York: ACM
- Hassani BK. 2021. Societal bias reinforcement through machine learning: a credit scoring perspective. *AI Ethics* 1(3):239–47
- Hertwig R, Engel C. 2016. Homo ignorans: deliberately choosing not to know. *Perspect. Psychol. Sci.* 11(3):359–72
- Hidalgo CA, Orghian D, Canals JA, De Almeida F, Martín N. 2021. *How Humans Judge Machines*. Cambridge, MA: MIT Press
- Hohenstein J, Kizilcec RF, DiFranzo D, Aghajari Z, Mieczkowski H, et al. 2023. Artificial intelligence in communication impacts language and social relationships. *Sci. Rep.* 13:5487
- Hong JW, Wang Y, Lanz P. 2020. Why is artificial intelligence blamed more? Analysis of faulting artificial intelligence for self-driving car accidents in experimental settings. *Int. J. Hum.-Comput. Interact.* 36(18):1768–74
- Hsieh TY, Cross ES. 2022. People’s dispositional cooperative tendencies towards robots are unaffected by robots’ negative emotional displays in prisoner’s dilemma games. *Cogn. Emot.* 36(5):995–1019
- Huang K, Greene JD, Bazerman M. 2019. Veil-of-ignorance reasoning favors the greater good. *PNAS* 116(48):23989–95
- Ishowo-Oloko F, Bonnefon JF, Soroye Z, Crandall J, Rahwan I, Rahwan T. 2019. Behavioural evidence for a transparency–efficiency tradeoff in human–machine cooperation. *Nat. Mach. Intell.* 1(11):517–21
- Jagatic TN, Johnson NA, Jakobsson M, Menczer F. 2007. Social phishing. *Commun. ACM* 50(10):94–100
- Jago AS, Laurin K. 2022. Assumptions about algorithms’ capacity for discrimination. *Pers. Soc. Psychol. Bull.* 48(4):582–95

- Jakesch M, French M, Ma X, Hancock JT, Naaman M. 2019. AI-mediated communication: how the perception that profile text was written by AI affects trustworthiness. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, Pap. 239. New York: ACM
- Kalra N, Groves DG. 2017. *The Enemy of Good: Estimating the Cost of Waiting for Nearly Perfect Automated Vehicles*. Santa Monica, CA: RAND
- Kalra N, Paddock SM. 2016. Driving to safety: How many miles of driving would it take to demonstrate autonomous vehicle reliability? *Transp. Res. A Policy Pract.* 94:182–93
- Karpus J, Krüger A, Verba JT, Bahrami B, Deroy O. 2021. Algorithm exploitation: Humans are keen to exploit benevolent AI. *iScience* 24(6):102679
- Kleinberg J, Lakkaraju H, Leskovec J, Ludwig J, Mullainathan S. 2018. Human decisions and machine predictions. *Q. J. Econ.* 133(1):237–93
- Kleinberg J, Mullainathan S, Raghavan M. 2017. *Inherent trade-offs in the fair determination of risk scores*. Paper presented at the 8th Innovations in Theoretical Computer Science Conference, Berkeley, CA, Jan. 9–11
- Köbis NC, Verschuere B, Bereby-Meyer Y, Rand D, Shalvi S. 2019. Intuitive honesty versus dishonesty: meta-analytic evidence. *Perspect. Psychol. Sci.* 14(5):778–96
- Komatsu T, Malle BF, Scheutz M. 2021. Blaming the reluctant robot: parallel blame judgments for robots in moral dilemmas across US and Japan. In *Proceedings of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, pp. 63–72. New York: ACM
- Kozyreva A, Herzog SM, Lewandowsky S, Hertwig R, Lorenz-Spreen P, et al. 2023. Resolving content moderation dilemmas between free speech and harmful misinformation. *PNAS* 120:e2210666120
- Krügel S, Uhl M. 2022. Autonomous vehicles and moral judgments under risk. *Transp. Res. A Policy Pract.* 155:1–10
- Lima G, Grgić-Hlača N, Cha M. 2021. Human perceptions on moral responsibility of AI: a case study in AI-assisted bail decision-making. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, Pap. 235. New York: ACM
- Liu P, Du M, Li T. 2021. Psychological consequences of legal responsibility misattribution associated with automated vehicles. *Ethics Inform. Technol.* 23(4):763–76
- Liu P, Du Y. 2022. Blame attribution asymmetry in human–automation cooperation. *Risk Anal.* 42(8):1769–83
- Liu P, Du Y, Xu Z. 2019a. Machines versus humans: people’s biased responses to traffic accidents involving self-driving vehicles. *Accid. Anal. Prevent.* 125:232–40
- Liu P, Liu J. 2021. Selfish or utilitarian automated vehicles? Deontological evaluation and public acceptance. *Int. J. Hum.–Comput. Interact.* 37(13):1231–42
- Liu P, Yang R, Xu Z. 2019b. How safe is safe enough for self-driving vehicles? *Risk Anal.* 39(2):315–25
- Liu Y, Mittal A, Yang D, Bruckman A. 2022. Will AI console me when I lose my pet? Understanding perceptions of AI-mediated email writing. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, Pap. 474. New York: ACM
- Longoni C, Bonezzi A, Morewedge CK. 2019. Resistance to medical artificial intelligence. *J. Consum. Res.* 46(4):629–50
- Longoni C, Cian L, Kyung EJ. 2022. Algorithmic transference: People overgeneralize failures of AI in the government. *J. Market. Res.* 60(1):170–88
- Luetge C. 2017. The German ethics code for automated and connected driving. *Philos. Technol.* 30(4):547–58
- Makovi K, Sargsyan A, Li W, Bonnefon JF, Rahwan T. 2023. Trust within human-machine collectives depends on the perceived consensus about cooperative norms. *Nat. Commun.* 14:3108
- Malle BF, Guglielmo S, Voiklis J, Monroe AE. 2022. Cognitive blame is socially shaped. *Curr. Dir. Psychol. Sci.* 31(2):169–76
- Malle BF, Scheutz M, Arnold T, Voiklis J, Cusimano C. 2015. Sacrifice one for the good of many? People apply different moral norms to human and robot agents. In *2015 10th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 117–24. Piscataway, NJ: IEEE
- March C. 2021. Strategic interactions between humans and artificial intelligence: lessons from experiments with computer players. *J. Econ. Psychol.* 87:102426
- Martin R, Kusev P, Van Schaik P. 2021. Autonomous vehicles: how perspective-taking accessibility alters moral judgments and consumer purchasing behavior. *Cognition* 212:104666

- Martinho A, Herber N, Kroesen M, Chorus C. 2021. Ethical issues in focus by the autonomous vehicles industry. *Transp. Rev.* 41(5):556–77
- Mathur MB, Reichling DB. 2016. Navigating a social world with robot partners: a quantitative cartography of the uncanny valley. *Cognition* 146:22–32
- Mayer MM, Bell R, Buchner A. 2021. Self-protective and self-sacrificing preferences of pedestrians and passengers in moral dilemmas involving autonomous vehicles. *PLOS ONE* 16(12):e0261673
- McAllister A. 2016. Stranger than science fiction: the rise of AI interrogation in the dawn of autonomous robots and the need for an additional protocol to the UN Convention Against Torture. *Minn. Law Rev.* 101:2527–73
- McCallum S, Vallance C. 2022. Start-up denies using tech to turn call centre accents “white.” *BBC News*, Aug. 26. <https://www.bbc.com/news/technology-62633188#>
- Miller T. 2019. Explanation in artificial intelligence: insights from the social sciences. *Artif. Intell.* 267:1–38
- Mittelstadt B. 2019. Principles alone cannot guarantee ethical AI. *Nat. Mach. Intell.* 1(11):501–7
- Moor JH. 2006. The nature, importance, and difficulty of machine ethics. *IEEE Intell. Syst.* 21(4):18–21
- Morley J, Floridi L, Kinsey L, Elhalal A. 2020. From what to how: an initial review of publicly available AI ethics tools, methods and research to translate principles into practices. *Sci. Eng. Ethics* 26(4):2141–68
- Mullainathan S. 2019. Biased algorithms are easier to fix than biased people. *New York Times*, Dec. 6. <https://www.nytimes.com/2019/12/06/business/algorithm-bias-fix.html>
- Nielsen YA, Pfattheicher S, Keijsers M. 2022a. Prosocial behavior toward machines. *Curr. Opin. Psychol.* 43:260–65
- Nielsen YA, Thielmann I, Zettler I, Pfattheicher S. 2022b. Sharing money with humans versus computers: on the role of honesty-humility and (non-)social preferences. *Soc. Psychol. Pers. Sci.* 13(6):1058–68
- Nightingale SJ, Farid H. 2022. AI-synthesized faces are indistinguishable from real faces and more trustworthy. *PNAS* 119(8):e2120481119
- Noy IY, Shinar D, Horrey WJ. 2018. Automated driving: safety blind spots. *Saf. Sci.* 102:68–78
- Nussberger AM, Luo L, Celis LE, Crockett MJ. 2022. Public attitudes value interpretability but prioritize accuracy in artificial intelligence. *Nat. Commun.* 13:5821
- O’Leary DE. 2019. Google’s duplex: pretending to be human. *Intell. Syst. Account. Finance Manag.* 26(1):46–53
- O’Neil C. 2017. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. New York: Crown
- Ostermaier A, Uhl M. 2017. Spot on for liars! How public scrutiny influences ethical behavior. *PLOS ONE* 12(7):e0181682
- Pammer K, Gauld C, McKerral A, Reeves C. 2021. “They have to be better than human drivers!” Motorcyclists’ and cyclists’ perceptions of autonomous vehicles. *Transp. Res. F Traffic Psychol. Behav.* 78:246–58
- Pammer K, Predojevic H, McKerral A. 2023. Humans vs. machines: Motorcyclists and car drivers differ in their opinion and trust of self-drive vehicles. *Transp. Res. F Traffic Psychol. Behav.* 92:143–54
- Pauketat JV, Anthis JR. 2022. Predicting the moral consideration of artificial intelligences. *Comput. Hum. Behav.* 136:107372
- Pleiss G, Raghavan M, Wu F, Kleinberg J, Weinberger KQ. 2017. On fairness and calibration. *Adv. Neural Inform. Process. Syst.* 30:5684–93
- Rauhut H. 2013. Beliefs about lying and spreading of dishonesty: undetected lies and their constructive and destructive social dynamics in dice experiments. *PLOS ONE* 8(11):e77878
- Rebitschek FG, Gigerenzer G, Wagner GG. 2021. People underestimate the errors made by algorithms for credit scoring and recidivism prediction but accept even fewer errors. *Sci. Rep.* 11:20171
- Rosenthal-von der Pütten AM, Krämer NC, Hoffmann L, Sobieraj S, Eimler SC. 2013. An experimental study on emotional reactions towards a robot. *Int. J. Soc. Robot.* 5(1):17–34
- Samuel S, Yalooz S, Yamani Y, Valluru K, Fisher DL. 2020. Ethical decision making behind the wheel—a driving simulator study. *Transp. Res. Interdiscip. Perspect.* 5:100147
- Sandoval EB, Brandstetter J, Obaid M, Bartneck C. 2016. Reciprocity in human-robot interaction: a quantitative approach through the prisoner’s dilemma and the ultimatum game. *Int. J. Soc. Robot.* 8(2):303–17

- Santoni de Sio F. 2021. The European Commission report on ethics of connected and automated vehicles and the future of ethics of transportation. *Ethics Inform. Technol.* 23(4):713–26
- Savulescu J, Gyngell C, Kahane G. 2021. Collective reflective equilibrium in practice (CREP) and controversial novel technologies. *Bioethics* 35(7):652–63
- Saxena NA, Huang K, DeFilippis E, Radanovic G, Parkes DC, Liu Y. 2020. How do fairness definitions fare? Testing public attitudes towards three algorithmic definitions of fairness in loan allocations. *Artif. Intell.* 283:103238
- Schwarting W, Pierson A, Alonso-Mora J, Karaman S, Rus D. 2019. Social behavior for autonomous vehicles. *PNAS* 116(50):24972–78
- Schwitzgebel E, Cushman F. 2015. Philosophers' biased judgments persist despite training, expertise and reflection. *Cognition* 141:127–37
- Seering J, Flores JP, Savage S, Hammer J. 2018. The social roles of bots: evaluating impact of bots on discussions in online communities. In *Proceedings of the ACM Human-Computer Interaction*, Vol. 2, Pap. 157. New York: ACM
- Seymour J, Tully P. 2016. *Weaponizing data science for social engineering: automated E2E spear phishing on Twitter*. Black Hat Video. <https://www.youtube.com/watch?v=-Y-xLQy0UuQ>
- Shank DB, DeSanti A, Maninger T. 2019. When are artificial intelligence versus human agents faulted for wrongdoing? Moral attributions after individual and joint decisions. *Inform. Commun. Soc.* 22(5):648–63
- Shao C, Ciampaglia GL, Varol O, Yang KC, Flammini A, Menczer F. 2018. The spread of low-credibility content by social bots. *Nat. Commun.* 9:4787
- Shariff A, Bonnefon JF, Rahwan I. 2017. Psychological roadblocks to the adoption of self-driving vehicles. *Nat. Hum. Behav.* 1(10):694–96
- Shariff A, Bonnefon JF, Rahwan I. 2021. How safe is safe enough? Psychological mechanisms underlying extreme safety demands for self-driving cars. *Transp. Res. C Emerg. Technol.* 126:103069
- Silver D, Huang A, Maddison CJ, Guez A, Sifre L, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 529(7587):484–89
- Smith A. 2019. *Public attitudes toward computer algorithms*. Rep., Pew Res. Cent., Washington, DC. <https://www.pewresearch.org/internet/2018/11/16/public-attitudes-toward-computer-algorithms>
- Srinivasan R, Sarial-Abi G. 2021. When algorithms fail: consumers' responses to brand harm crises caused by algorithm errors. *J. Market.* 85(5):74–91
- Srivastava M, Heidari H, Krause A. 2019. Mathematical notions versus human perception of fairness: a descriptive approach to fairness for machine learning. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2459–68. New York: ACM
- Starke C, Baleis J, Keller B, Marcinkowski F. 2022. Fairness perceptions of algorithmic decision-making: a systematic review of the empirical literature. *Big Data Soc.* 9(2):20539517221115189
- Stella M, Ferrara E, De Domenico M. 2018. Bots increase exposure to negative and inflammatory content in online social systems. *PNAS* 115(49):12435–40
- Suzuki Y, Galli L, Ikeda A, Itakura S, Kitazaki M. 2015. Measuring empathy for human and robot hand pain using electroencephalography. *Sci. Rep.* 5:15924
- Takaguchi K, Kappes A, Yearsley JM, Sawai T, Wilkinson DJ, Savulescu J. 2022. Personal ethical settings for driverless cars and the utility paradox: an ethical analysis of public attitudes in UK and Japan. *PLOS ONE* 17(11):e0275812
- Thomas PS, Castro da Silva B, Barto AG, Giguere S, Brun Y, Brunskill E. 2019. Preventing undesirable behavior of intelligent machines. *Science* 366(6468):999–1004
- Tsvetkova M, García-Gavilanes R, Floridi L, Yasseri T. 2017. Even good bots fight: the case of Wikipedia. *PLOS ONE* 12(2):e0171774
- Van Zant AB, Kray LJ. 2014. “I can't lie to your face”: Minimal face-to-face interaction promotes honesty. *J. Exp. Soc. Psychol.* 55:234–38
- Villani V, Pini F, Leali F, Secchi C. 2018. Survey on human-robot collaboration in industrial settings: safety, intuitive interfaces and applications. *Mechatronics* 55:248–66
- von Schenk A, Klockmann V, Köbis N. 2022. *Social preferences towards machines and humans*. Work. Pap., Max Planck Inst. Hum. Dev., Berlin, Ger.

- Wegner D, Gray K. 2016. *The Mind Club: Who Thinks, What Feels, and Why It Matters*. New York: Viking
- Wellman MP, Rajan U. 2017. Ethical issues for autonomous trading agents. *Minds Mach.* 27:609–24
- Wotton ME, Bennett JM, Modesto O, Challinor KL, Prabhakaran P. 2022. Attention all “drivers”: You could be to blame, no matter your behaviour or the level of vehicle automation. *Transp. Res. F Traffic Psychol. Behav.* 87:219–35
- Zhu A, Yang S, Chen Y, Xing C. 2022. A moral decision-making study of autonomous vehicles: Expertise predicts a preference for algorithms in dilemmas. *Pers. Individ. Differ.* 186:111356
- Zlotowski J, Sumioka H, Nishio S, Glas DF, Bartneck C, Ishiguro H. 2016. Appearance of a robot affects the impact of its behaviour on perceived trustworthiness and empathy. *Paladyn J. Behav. Robot.* 7:55–66



# Contents

The Neurobiology of Activational Aspects of Motivation: Exertion of Effort, Effort-Based Decision Making, and the Role of Dopamine <i>John D. Salamone and Mercè Correa</i> .....	1
Sexual Incentive Motivation and Sexual Behavior: The Role of Consent <i>Anders Ågmo and Ellen Laan</i> .....	33
A Systematic Review of Implementation Research on Determinants and Strategies of Effective HIV Interventions for Men Who Have Sex with Men in the United States <i>Brian Mustanski, Artur Queiroz, James L. Merle, alithia zamantakis, Juan Pablo Zapata, Dennis H. Li, Nanette Benbow, Maria Pyra, and Justin D. Smith</i> .....	55
Music Training and Nonmusical Abilities <i>E. Glenn Schellenberg and César F. Lima</i> .....	87
Serial Dependence in Perception <i>Guido Marco Cicchini, Kyriaki Mikellidou, and David Charles Burr</i> .....	129
What Does the Human Olfactory System Do, and How Does It Do It? <i>Gülce Nazlı Dikeçligil and Jay A. Gottfried</i> .....	155
The Relation Between Attention and Memory <i>Nelson Cowan, Chenye Bao, Brittney M. Bishop-Chrzanowski, Amy N. Costa, Nathaniel R. Greene, Dominic Guitard, Chenyuan Li, Madison L. Musich, and Zebra E. Ünal</i> .....	183
Modeling Similarity and Psychological Space <i>Brett D. Roads and Bradley C. Love</i> .....	215
Metacognition and Confidence: A Review and Synthesis <i>Stephen M. Fleming</i> .....	241
Beyond the Tricks: The Science and Comparative Cognition of Magic <i>Elias Garcia-Pelegrin, Alexandra K. Schnell, Clive Wilkins, and Nicola S. Clayton</i> ....	269
Moral Improvement of Self, Social Relations, and Society <i>Colin Wayne Leach and Arti Iyer</i> .....	295



Social Media and Morality <i>Jay J. Van Bavel, Claire E. Robertson, Kareena del Rosario, Jesper Rasmussen, and Steve Rathje</i> .....	311
Norm Dynamics: Interdisciplinary Perspectives on Social Norm Emergence, Persistence, and Change <i>Michele J. Gelfand, Sergey Gavrilets, and Nathan Nunn</i> .....	341
Pursuing Safety in Social Connection: A Flexibly Fluid Perspective on Risk Regulation in Relationships <i>Sandra L. Murray and Gabriela S. Pascuzzi</i> .....	379
Knowledge Transfer Within Organizations: Mechanisms, Motivation, and Consideration <i>Linda Argote</i> .....	405
The Neuroscience of Human and Artificial Intelligence Presence <i>Lasana T. Harris</i> .....	433
How Can People Become Happier? A Systematic Review of Preregistered Experiments <i>Dunigan Folk and Elizabeth Dunn</i> .....	467
Cultural Psychology: Beyond East and West <i>Shinobu Kitayama and Cristina E. Salvador</i> .....	495
Achievement Goals: A Social Influence Cycle <i>Fabrizio Butera, Benoît Dompnier, and Céline Darnon</i> .....	527
Why We Should Stop Trying to Fix Women: How Context Shapes and Constrains Women's Career Trajectories <i>Michelle K. Ryan and Thekla Morgenroth</i> .....	555
Resilience and Disaster: Flexible Adaptation in the Face of Uncertain Threat <i>George A. Bonanno, Shuquan Chen, Robini Bagrodia, and Isaac R. Galatzer-Levy</i> .....	573
Psychological Flexibility, Chronic Pain, and Health <i>Lance M. McCracken</i> .....	601
Computational Social Psychology <i>Fiery Cushman</i> .....	625
The Moral Psychology of Artificial Intelligence <i>Jean-François Bonnefon, Iyad Rahwan, and Azim Shariff</i> .....	653

## Indexes

Cumulative Index of Contributing Authors, Volumes 65–75 .....	677
Cumulative Index of Article Titles, Volumes 65–75 .....	682

## Errata

An online log of corrections to *Annual Review of Psychology* articles may be found at <http://www.annualreviews.org/errata/psych>

## Related Articles

From the *Annual Review of Clinical Psychology*, Volume 19 (2023)

A Clinical Psychologist Who Studies Alcohol

*Kenneth J. Sher*

Community Mental Health Services for American Indians and Alaska Natives:  
Reconciling Evidence-Based Practice and Alter-Native Psy-ence

*Joseph P. Gone*

Culturally Responsive Cognitive Behavioral Therapy for Ethnically Diverse  
Populations

*Stanley J. Huey Jr., Alayna L. Park, Chardée A. Galán, and Crystal X. Wang*

What Four Decades of Meta-Analysis Have Taught Us About Youth  
Psychotherapy and the Science of Research Synthesis

*John R. Weisz, Katherine E. Ventura-Conerly, Olivia M. Fitzpatrick,  
Jennifer A. Frederick, and Mei Yi Ng*

Evaluation of Pressing Issues in Ecological Momentary Assessment

*Arthur A. Stone, Stefan Schneider, and Joshua M. Smyth*

Machine Learning and the Digital Measurement of Psychological Health

*Isaac R. Galatzer-Levy and Jukka-Pekka Onnela*

The Questionable Practice of Partialing to Refine Scores on and Inferences About  
Measures of Psychological Constructs

*Rick H. Hoyle, Donald R. Lynam, Joshua D. Miller, and Jolynn Pek*

Eating Disorders in Boys and Men

*Tiffany A. Brown and Pamela K. Keel*

Mental Health of Transgender and Gender Diverse Youth

*Natalie M. Wittlin, Laura E. Kuper, and Kristina R. Olson*

Behavioral Interventions for Children and Adults with Tic Disorder

*Douglas W. Woods, Michael B. Himle, Jordan T. Stiede, and Brandon X. Pitts*

The Garrett Lee Smith Memorial Act: A Description and Review of the Suicide  
Prevention Initiative

*David B. Goldston and Christine Walrath*

Racism and Social Determinants of Psychosis

*Deidre M. Anglin*

Developmental Consequences of Intimate Partner Violence on Children

*G. Anne Bogat, Alytia A. Levendosky, and Kara Cochran*

Psychoneuroimmunology: An Introduction to Immune-to-Brain Communication and Its Implications for Clinical Psychology

*Julienne E. Bower and Kate R. Kublman*

Racial, Ethnic, and Cultural Resilience Factors in African American Youth Mental Health

*Enrique W. Neblett Jr.*

Acculturation and Psychopathology

*Gail M. Ferguson, José M. Causadias, and Tori S. Simenec*

Posttraumatic Stress Disorder in Refugees

*Richard A. Bryant, Angela Nickerson, Naser Morina, and Belinda Liddell*

Risk and Resilience Among Children with Incarcerated Parents: A Review and Critical Reframing

*Elizabeth I. Johnson and Joyce A. Arditti*

Supernatural Attributions: Seeing God, the Devil, Demons, Spirits, Fate, and Karma as Causes of Events

*Julie J. Exline and Joshua A. Wilt*

From the *Annual Review of Developmental Psychology*, Volume 5 (2023)

Navigating an Unforeseen Pathway

*Margaret Beale Spencer*

Prenatal Substance Exposure

*Rina D. Eiden, Kristin J. Perry, Miglena Y. Ivanova, and Rachel C. Marcus*

Neurodevelopment of Attention, Learning, and Memory Systems in Infancy

*Tess Allegra Forest and Dima Amso*

The Representation of Third-Party Helping Interactions in Infancy

*Laura Schlingloff-Nemecz, Denis Tatone, and Gergely Csibra*

A Developmental Social Neuroscience Perspective on Infant Autism Interventions

*Geraldine Dawson, Amber D. Rieder, and Mark H. Johnson*

Intervening Early: Socioemotional Interventions Targeting the Parent–Infant Relationship

*Mary Dozier and Kristin Bernard*

Growing Up, Learning Race: An Integration of Research on Cognitive Mechanisms and Socialization in Context

*Diane Hughes, Blair Cox, and Sobini Das*

Social Identities and Intersectionality: A Conversation About the What and the How of Development

*Margarita Azmitia, Paulette D. Garcia Peraza, and Saskias Casanova*

Children's Acquisition and Application of Norms

*Marco F.H. Schmidt and Hannes Rakoczy*

A Rational Account of Cognitive Control Development in Childhood

*Nikolaus Steinbeis*

A Two-Hit Model of Behavioral Inhibition and Anxiety

*Brendan Oslund and Koraly Pérez-Edgar*

Developmental Neuroimaging of Cognitive Flexibility: Update and Future Directions

*Lauren B. Kupis and Lucina Q. Uddin*

A Neuroecosocial Perspective on Adolescent Development

*Suparna Choudhury, Blanca Piera Pi-Sunyer, and Sarah-Jayne Blakemore*

Poverty, Brain Development, and Mental Health: Progress, Challenges, and Paths Forward

*Christopher S. Monk and Felicia A. Hardi*

The Study of Early Child Care and Youth Development (SECCYD): Studying Development from Infancy to Adulthood

*Deborah Lowe Vandell and Zebra Gülseven*

Bridging the Divide: Tackling Tensions Between Life-Course Epidemiology and Causal Inference

*Gabriel L. Schwartz and M. Maria Glymour*

The Functioning of Offspring of Depressed Parents: Current Status, Unresolved Issues, and Future Directions

*Ian H. Gotlib, Jessica L. Buttmann, and Jonas G. Miller*

Emotion Regulation in Couples Across Adulthood

*Claudia M. Haase*

From the *Annual Review of Neuroscience*, Volume 46 (2023)

Therapeutic Potential of PTBP1 Inhibition, If Any, Is Not Attributed to Glia-to-Neuron Conversion

*Lei-Lei Wang and Chun-Li Zhang*

How Flies See Motion

*Alexander Borst and Lukas N. Groschner*

Meningeal Mechanisms and the Migraine Connection

*Dan Levy and Michael A. Moskowitz*

Cholesterol Metabolism in Aging and Age-Related Disorders

*Gesine Saher*

Spinal Interneurons: Diversity and Connectivity in Motor Control

*Mohini Sengupta and Martha W. Bagnall*

Astrocyte Endfeet in Brain Function and Pathology: Open Questions

*Blanca Díaz-Castro, Stefanie Robel, and Anusha Misbra*

- Circadian Rhythms and Astrocytes: The Good, the Bad, and the Ugly  
*Michael H. Hastings, Marco Brancaccio, Maria F. Gonzalez-Aponte,  
and Erik D. Herzog*
- Therapeutic Potential of PTB Inhibition Through Converting Glial Cells  
to Neurons in the Brain  
*Xiang-Dong Fu and William C. Mobley*
- How Instructions, Learning, and Expectations Shape Pain and Neurobiological  
Responses  
*Lauren Y. Atlas*
- Cognition from the Body-Brain Partnership: Exaptation of Memory  
*György Buzsáki and David Tingley*
- Neural Circuits for Emotion  
*Meryl Malezieux, Alexandra S. Klein, and Nadine Gogolla*
- The Computational and Neural Bases of Context-Dependent Learning  
*James B. Heald, Daniel M. Wolpert, and Máté Lengyel*
- Integration of Feedforward and Feedback Information Streams in the Modular  
Architecture of Mouse Visual Cortex  
*Andreas Burkhalter, Rinaldo D. D'Souza, Weiqing Ji, and Andrew M. Meier*
- How Do You Build a Cognitive Map? The Development of Circuits  
and Computations for the Representation of Space in the Brain  
*Flavio Donato, Anja Xu Schwartzlose, and Renan Augusto Viana Mendes*
- Cortical Integration of Vestibular and Visual Cues for Navigation, Visual  
Processing, and Perception  
*Sepideh Keshavarzi, Mateo Velez-Fort, and Troy W. Margrie*
- Neural Control of Sexually Dimorphic Social Behavior: Connecting Development  
to Adulthood  
*Margaret M. McCarthy*
- Deep Brain Stimulation for Obsessive-Compulsive Disorder and Depression  
*Sameer A. Sheth and Helen S. Mayberg*
- Striosomes and Matrisomes: Scaffolds for Dynamic Coupling of Volition  
and Action  
*Ann M. Graybiel and Ayano Matsushima*
- Specialized Networks for Social Cognition in the Primate Brain  
*Ben Deen, Caspar M. Schwiedrzik, Julia Sliwa, and Winrich A. Freiwald*
- Neural Networks for Navigation: From Connections to Computations  
*Rachel I. Wilson*

From the *Annual Review of Organizational Psychology and Organizational Behavior*,  
Volume 10 (2023)

Changes in Perspective and Perspectives on Change: Reflections on a Career  
*Timothy A. Judge*

- Job Demands–Resources Theory: Ten Years Later  
*Arnold B. Bakker, Evangelia Demerouti, and Ana Sanz-Vergel*
- Psychological Safety Comes of Age: Observed Themes in an Established Literature  
*Amy C. Edmondson and Derrick P. Bransby*
- Employee Voice and Silence: Taking Stock a Decade Later  
*Elizabeth Wolfe Morrison*
- Understanding the Dynamic Interplay Between Actor and Context for Creativity: Progress and Desirable Directions  
*Jing Zhou and Inga J. Hoever*
- The Psychology of Entrepreneurship: Action and Process  
*Michael Frese and Michael M. Gielnik*
- Laying the Foundation for the Challenge–Hindrance Stressor Framework 2.0  
*Nathan P. Podsakoff, Kristen J. Freiburger, Philip M. Podsakoff, and Christopher C. Rosen*
- Crisis Leadership  
*Ronald E. Riggio and Toby Newstead*
- Meta-Analysis in Organizational Research: A Guide to Methodological Options  
*Scott B. Morris*
- Developing Self-Awareness: Learning Processes for Self- and Interpersonal Growth  
*Manuel London, Valerie I. Sessa, and Loren A. Shelley*
- Understanding Decent Work and Meaningful Work  
*David L. Blustein, Evgenia I. Lysova, and Ryan D. Duffy*
- Innovations in Sampling: Improving the Appropriateness and Quality of Samples in Organizational Research  
*Michael J. Zickar and Melissa G. Keith*
- Leading Virtually  
*Bradford S. Bell, Kristie L. McAlpine, and N. Sharon Hill*
- Mental Health in the Workplace  
*E. Kevin Kelloway, Jennifer K. Dimoff, and Stephanie Gilbert*
- Is Justice Colorblind? A Review of Workplace Racioethnic Differences Through the Lens of Organizational Justice  
*Derek R. Avery, Alison V. Hall, McKenzie Preston, Enrica N. Ruggs, and Ella Washington*
- Leader Thinking, Follower Thinking: Leader Impacts on Follower Creative Performance  
*Michael D. Mumford, Mark Fichtel, Samantha England, and Tanner R. Newbold*
- Self-Reflection at Work: Why It Matters and How to Harness Its Potential and Avoid Its Pitfalls  
*Ethan Kross, Madeline Ong, and Ozlem Ayduk*

Employee Green Behavior as the Core of Environmentally Sustainable Organizations

*Hannes Zacher, Cort W. Rudolph, and Ian M. Katz*

Structural Equation Modeling in Organizational Research: The State of Our Science and Some Proposals for Its Future

*Michael J. Zyphur, Cavan V. Bonner, and Louis Tay*

Improving Workplace Judgments by Reducing Noise: Lessons Learned from a Century of Selection Research

*Scott Highhouse and Margaret E. Brooks*

From the *Annual Review of Public Health*, Volume 44 (2023)

A Literature Review of the Effects of Air Pollution on COVID-19 Health Outcomes Worldwide: Statistical Challenges and Data Visualization

*A. Bhaskar, J. Chandra, H. Hashemi, K. Butler, L. Bennett, Jacqueline Cellini, Danielle Braun, and Francesca Dominici*

On-the-Go Adaptation of Implementation Approaches and Strategies in Health: Emerging Perspectives and Research Opportunities

*Elvin H. Geng, Aaloke Mody, and Byron J. Powell*

Enhancing Capacity for Food and Nutrient Intake Assessment in Population Sciences Research

*Marian L. Neubouser, Ross L. Prentice, Lesley F. Tinker, and Johanna W. Lampe*

Innovations in Public Health Surveillance for Emerging Infections

*Peng Jia, Shiyong Liu, and Shujuan Yang*

Early Childhood Education: Health, Equity, and Economics

*Robert A. Hahn and W. Steven Barnett*

Environmental Justice: Where It Has Been, and Where It Might Be Going

*Merlin Chowkwanyun*

Health Misinformation Exposure and Health Disparities: Observations and Opportunities

*Brian G. Southwell, Jessica Otero Machuca, Sabrina T. Cherry, Melissa Burnside, and Nadine J. Barrett*

Leveraging Mobile Technology for Public Health Promotion: A Multidisciplinary Perspective

*Jennifer L. Hicks, Melissa A. Boswell, Tim Althoff, Alia J. Crum, Joy P. Ku, James A. Landay, Paula M.L. Moya, Elizabeth L. Murnane, Michael P. Snyder, Abby C. King, and Scott L. Delp*

When Moving Is the Only Option: The Role of Necessity Versus Choice for Understanding and Promoting Physical Activity in Low- and Middle-Income Countries

*Deborah Salvo, Alejandra Jáuregui, Deepti Adlakha, Olga L. Sarmiento, and Rodrigo S. Reis*



Climatic and Environmental Change, Migration, and Health

*Celia McMichael*

Promoting Health Equity Through Preventing or Mitigating the Effects of Gentrification: A Theoretical and Methodological Guide

*Helen V.S. Cole, Isabelle Anguelovski, Margarita Triguero-Mas, Roshanak Mehdipanah, and Mariana Arcaya*

Public Health Implications of Drought in a Climate Change Context: A Critical Review

*Coral Salvador, Raquel Nieto, Sergio M. Vicente-Serrano, Ricardo García-Herrera, Luis Gimeno, and Ana M. Vicedo-Cabrera*

Review of the Impact of Housing Quality on Inequalities in Health and Well-Being

*Philippa Howden-Chapman, Julie Bennett, Richard Edwards, David Jacobs, Kim Nathan, and David Ormandy*

Sustainable and Resilient Health Care in the Face of a Changing Climate

*Jodi D. Sherman, Andrea J. MacNeill, Paul D. Biddinger, Ozlem Ergun, Renee N. Salas, and Matthew J. Eckelman*

Cancers Attributable to Modifiable Risk Factors: A Road Map for Prevention

*Giulia Collatuzzo and Paolo Boffetta*

Public Health Preparedness for Extreme Heat Events

*Jeremy J. Hess, Nicole A. Errett, Glenn McGregor, Tania Busch Isaksen, Zachary S. Wettstein, Stefan K. Wheat, and Kristie L. Ebi*

The State of the US Public Health Workforce: Ongoing Challenges and Future Directions

*Jonathon P. Leider, Valerie A. Yeager, Chelsey Kirkland, Heather Krasna, Rachel Hare Bork, and Beth Resnick*

The Value and Impacts of Academic Public Health Departments

*Paul C. Erwin, Julie H. Grubaugh, Stephanie Mazzucca-Ragan, and Ross C. Brownson*

Community Health Worker Integration with and Effectiveness in Health Care and Public Health in the United States

*Molly Knowles, Aidan P. Crowley, Aditi Vasan, and Shreya Kangovi*

Multilevel Determinants of Digital Health Equity: A Literature Synthesis to Advance the Field

*Courtney R. Lyles, Oanh Kieu Nguyen, Elaine C. Khoong, Adrian Aguilera, and Urmimala Sarkar*

Public Health and Prisons: Priorities in the Age of Mass Incarceration

*David H. Cloud, Ilana R. Garcia-Grossman, Andrea Armstrong, and Brie Williams*

The Impacts of Paid Family and Medical Leave on Worker Health, Family Well-Being, and Employer Outcomes

*Ann Bartel, Maya Rossin-Slater, Christopher Rubm, Meredith Slopen, and Jane Waldfogel*

Using Rapid Randomized Trials to Improve Health Care Systems  
*Leora I. Horwitz and Holly A. Krelle*

From the *Annual Review of Vision Science*, Volume 9 (2023)

Envisioning a Woman Scientist  
*Suzanne P. McKee*

Disparities in Eye Care Access and Utilization: A Narrative Review  
*Joana E. Andoh, Agnes C. Ezekwesili, Kristen Nwanyanwu, and Angela Elam*

Pathophysiology of Retinopathy of Prematurity  
*M. Elizabeth Hartnett*

Emerging Pathogenic Viral Infections of the Eye  
*Ekta Rishi, Joanne Thomas, Tolulope Fashina, Lucas Kim, and Steven Yeh*

Visual Dysfunction in Diabetes  
*Erika D. Eggers*

Impact of Apps as Assistive Devices for Visually Impaired Persons  
*Shrinivas Pundlik, Prerana Shivshanker, and Gang Luo*

Suppressing Retinal Remodeling to Mitigate Vision Loss in Photoreceptor  
Degenerative Disorders  
*Richard H. Kramer*

Factors Affecting Stem Cell–Based Regenerative Approaches in Retinal  
Degeneration  
*Sachin H. Patel and Deepak A. Lamba*

Structure, Function, and Molecular Landscapes of the Aging Retina  
*Jeffrey D. Zhu, Sharma Pooja Tarachand, Qudrat Abdulwahab,  
and Melanie A. Samuel*

What Is a Preferred Retinal Locus?  
*Luminita Tarita-Nistor, Irina Sverdlichenko, and Mark S. Mandelcorn*

Eye Morphogenesis in Vertebrates  
*Macaulie A. Casey, Sarah Lusk, and Kristen M. Kwan*

An Expanding Role for Nonvisual Opsins in Extraocular Light Sensing Physiology  
*Mutabar Andrabi, Brian A. Upton, Richard A. Lang, and Shruti Vemaraju*

Predicting Visual Fixations  
*Matthias Kümmerer and Matthias Bethge*

Two Sides of the Same Coin: Efficient and Predictive Neural Coding  
*Michael B. Manookin and Fred Rieke*

Visual Representations: Insights from Neural Decoding  
*Amanda K. Robinson, Genevieve L. Quek, and Thomas A. Carlson*

Neuronal Representations Supporting Three-Dimensional Vision in Nonhuman  
Primates  
*Ari Rosenberg, Lowell W. Thompson, Raymond Doudlah, and Ting-Yu Chang*

Visual Functions of the Primate Superior Colliculus

*Ziad M. Hafed, Klaus-Peter Hoffmann, Chib-Yang Chen, and Amarendra R. Bogadhi*

Contributions of the Basal Ganglia to Visual Perceptual Decisions

*Long Ding*

Perception and Memory in the Ventral Visual Stream and Medial Temporal Lobe

*Chris B. Martin and Morgan D. Barense*

Using Natural Scenes to Enhance our Understanding of the Cerebral Cortex's  
Role in Visual Search

*Mark A. Segraves*

The Perceptual Science of Augmented Reality

*Emily A. Cooper*

Ultra-High Field Imaging of Human Visual Cognition

*Ke Jia, Rainer Goebel, and Zoe Kourtzi*

Are Deep Neural Networks Adequate Behavioral Models of Human Visual  
Perception?

*Felix A. Wichmann and Robert Geirbos*