



AI and the expert; a blueprint for the ethical use of opaque AI

Amber Ross¹

Received: 5 December 2021 / Accepted: 13 September 2022
© The Author(s), under exclusive licence to Springer-Verlag London Ltd., part of Springer Nature 2022

Abstract

The increasing demand for transparency in AI has recently come under scrutiny. The question is often posted in terms of “epistemic double standards”, and whether the standards for transparency in AI ought to be higher than, or equivalent to, our standards for ordinary human reasoners. I agree that the push for increased transparency in AI deserves closer examination, and that comparing these standards to our standards of transparency for other opaque systems is an appropriate starting point. I suggest that a more fruitful exploration of this question will involve a different comparison class. We routinely treat judgments made by highly trained experts in specialized fields as fair or well grounded even though—by the nature of expert/layperson division of epistemic labor—an expert will not be able to provide an explanation of the reasoning behind these judgments that makes sense to most other people. Regardless, laypeople are thought to be acting reasonably—and ethically—in deferring to the judgments of experts that concern their areas of specialization. I suggest that we reframe our question regarding the appropriate standards of transparency in AI as one that asks when, why, and to what degree it would be ethical to accept opacity in AI. I argue that our epistemic relation to certain opaque AI technology may be relevantly similar to the layperson’s epistemic relation to the expert in certain respects, such that the successful expert/layperson division of epistemic labor can serve as a blueprint for the ethical use of opaque AI.

Keywords AI Ethics · Opacity · Transparency · Explicability · Social epistemology · Expert testimony

1 Introduction

Does the widespread demand for increased transparency in AI impose an epistemic double standard on the judgments made by AI models? And if so, are those double standards justified? Should we hold AI technology to the same standards of transparency that we hold an ordinary human reasoner? These questions are beginning to receive attention in the AI ethics literature, but to date there is minimal consensus. Zerilli et al. (2019) argue that much of our current proposed regulations would hold AI to higher than normal—and higher than necessary—standards of transparency. Günther & Kasirzadeh (2022) hold that, while there may be a double standard for ordinary human judgments and judgments made by AI, those heightened standards for AI technology are appropriate.

Though they disagree on what the standards for AI transparency ought to be, all parties seem to accept that the standards to which they should be compared are our standards for transparency in the judgments of ordinary human reasoners. This makes sense, insofar as one’s own decision-making process is thought to be transparent to oneself, while the reasoning of other minds is notoriously opaque. And in high-stakes decisions, or contexts in which fairness is an issue, we certainly require at least some degree of explanation or transparency before we will accept a person’s judgment as fair and well grounded. Though we may not demand a full accounting of the reasoning process that ordinary people engage in when they make these judgments, our standards require that, at minimum, they ought to be able to provide an explanation of their reasoning that makes sense to most other people.

While I agree that the widespread push for increased transparency in AI deserves closer examination and that comparing these to our standards of transparency for other opaque systems is an appropriate starting point, I believe that a more fruitful exploration of this question will involve a different comparison class. While our most ubiquitous

A1 Amber Ross
A2 amber.ross@ufl.edu

A3 ¹ Department of Philosophy, The University of Florida,
A4 Gainesville, FL, USA

standards of transparency are those that apply to ordinary human reasoners making ordinary decisions, there is another familiar class of judgments to which these ordinary standards of transparency do not apply. We routinely treat judgments made by highly trained experts in specialized fields as legitimate or well-grounded even though—by the nature of expert/layperson division of epistemic labor—an expert will not be able to provide an explanation of the reasoning behind these judgments that makes sense to most other people. Despite this fact, most other people (those who are not experts in the particular specialized field) would be acting reasonably—and ethically—in deferring to the judgment of experts regarding matters that concern their area of specialization. I suggest that we might make progress on questions regarding the appropriate standards of transparency in AI by reframing the question as one that asks when, why, and to what degree it would be ethical to accept opacity in AI. As I will argue, our relation to some opaque AI technology may be sufficiently similar to the ordinary layperson's relation to the specialized expert such that analyzing the successful expert/layperson epistemic relation may provide us with a blueprint for how to best utilize opaque AI, both practically and ethically.

The general organization of this paper will be as follows: In Sect. 2, I will discuss the general value of allowing for the kind of opacity that exists in the expert/layperson relation. In Sect. 3, I will address the value of transparency in decision-making, focusing on automated decision makers (ADMs) and the problem of bias in machine learning. In Sect. 4 I will explore areas of ethical concern *beyond* bias. Fairness is one value among many that must be considered when developing guidelines for the ethical use of AI. I believe an overly concentrated focus on the problem of bias in AI has drawn our attention away from other values that need to be considered in a full-cost accounting of our use of AI. It is the presence of these additional considerations that show why, in certain cases, allowing for opacity in AI technology may be ethically preferable to a constant pursuit of transparency. In Sect. 5, I will argue that the call for transparency in AI is mainly in service of a separate end—that transparency serves as a proxy for the trustworthiness of opaque processes, and increasing transparency aims at establishing appropriate levels of trust between stakeholders and opaque AI. If this is correct, we may be ethically permitted to utilize opaque AI technology provided that this trust and trustworthiness can be established through alternate means. In Sect. 6, I will give an overview of several fundamental features of the expert/layperson relation and make a case for the possibility that the relation between stakeholders and opaque AI could display these features as well. These features will provide a skeletal blueprint for the ethical use of opaque AI. In Sect. 7, I will suggest preliminary guidelines for evaluating contexts in which it may be ethical to employ opaque AI

technology, consistent with the blueprint adapted from the successful expert/layperson relation.

2 The value of harnessing opaque processes

As a society, we reap enormous benefits from relying on—or deferring to—expert judgments, especially in high-stakes contexts. Our division of epistemic labor allows laypeople to benefit from the knowledge and judgments of specialized experts without understanding *how* the experts arrive at these judgments nor *why* those judgments are justified. If there is such a thing as scientific progress, discovering how to effectively utilize this division of epistemic labor is the foundation of that progress.

Our reliance on opaque expert reasoning is so common that it usually passes without our notice. It may be as trivial as relying on the weather forecast when planning a vacation, or as significant as deciding whether to evacuate our homes (risking our lives and livelihoods) because we know we are in the path of a hurricane. In modern society, one does not need to understand the nature of carbon monoxide or nuclear reactions to know that certain levels of CO in the home can be deadly, or that certain nuclear power plants are safe to live near. We can make ethically responsible decisions, including high-stakes decisions, without fully understanding the reasoning process on which we are basing our decision, because it is both epistemically and ethically responsible for us to defer to experts in these matters.

For the vast majority of society, the evidence and reasoning processes of any expert in a specialized field is opaque, a genuine “black box”. Though it is often in our best interest to defer to these experts' judgments, in doing so we are accepting the outcome of a process that we are aware we do not understand. We individuals who are not experts in a particular specialized field—can know far more than we have the capacity to understand, because relying on expert opinion is a reliable way to build knowledge and an ethically responsible way to decide how to act. A medical expert can only make their reasoning and evidence understandable to a layperson to a certain degree; for that reasoning to be fully transparent to the patient, the patient would need to undergo training similar to that which the doctor underwent to become an expert in their field. This is obviously impractical and undesirable. Instead, we routinely rely on reasoning that we do not understand—especially in high-stakes situations—and this practice is indispensable to modern life. We defer to the judgments of medical doctors, structural engineers, epidemiologists, meteorologists, and computer scientists on a daily basis, and we do so precisely because we know we do not know what qualifies as good evidence or good reasoning in these highly specialized fields.

160 Just as human expertise is most useful in areas where
 161 sound judgments require extended and complex training
 162 in specialized fields (making the required reasoning
 163 opaque to most), AI is most useful in areas where its
 164 speed and capacity for data processing greatly surpasses
 165 human abilities—the same factors that make certain AI
 166 technology opaque. And just as the judgment of experts
 167 is most valuable in high-stakes situations, the maximal
 168 benefit we can derive from AI will be in its application to
 169 areas that are central to human welfare (areas such as
 170 health, agriculture, climate, and public safety). The power
 171 of AI is a double-edged sword. Its extraordinary speed and
 172 unconventional data-processing methods are the same
 173 factors that can make the most powerful AI opaque to its
 174 users and stakeholders, creating ethical concerns
 175 regarding whether it ought to be used in the very areas in
 176 which it could potentially provide the most benefit. The
 177 more knowledge we are ethically required to have
 178 regarding how AI technology works when it operates in
 179 a particular domain, the less likely it is that we will be
 ethically permitted to use AI in that domain.

3 Opacity and the problem of algorithmic bias

180 The call for transparency in AI aims at safeguarding and
 181 improving human welfare—in particular, by
 182 protecting vulnerable groups who are most often
 183 harmed by opaque AI technology and marginalized in AI
 184 development. This goal is and should be a top priority in
 185 AI regulation. The speed and processing power of AI not
 186 only comes at an episodic cost; as we have learned, our
 187 limited epistemic access to certain AI models can bring
 188 with it ethical costs as well. In 2016, investigative
 189 journalists at ProPublica published an article that
 190 exposed apparent racial bias in the popular risk-
 191 assessment software COMPAS, used to aid judicial
 decision-making regarding individuals' risk of recidivism
 and eligibility for parole. In 2018, Reuters¹ revealed that
 the AI hiring algorithm in development at Google
 showed a strong gender bias.

192 The push to integrate these ADMs into areas such as
 193 recidivism risk assessment, loan approval, and hiring
 194 practices, has exposed a tension between two
 195 worthwhile goals:

- 196 (i) increased efficiency in important decision-making
 197 processes and (ii) protecting individuals' rights by
 198 ensuring such decisions are based only upon ethically
 199 appropriate considerations. This tension can become
 200 more problem-

204 atic when the AI technology involved is opaque—when the

205 methods by which the AI arrives at a decision cannot be
 206 tracked by the relevant parties, whether AI practitioner
 207 or stakeholder.

208 The most powerful AI models—such as deep learning
 209 models and models involving vast parameters—are also
 210 the least comprehensible. While the engineers involved
 211 in creating ADMs like COMPAS may be aware of the
 212 content of the training dataset and the parameters at the
 213 time of use, the precise role these play in generating the
 214 ADM's output often remains unknown. For very
 215 complex models, there may be no one (neither AI
 216 practitioner nor stakeholder) who understands the
 217 actual relevance of each datum to the ADM's eventual
 218 prediction. As Riberio (2016) writes, "... if hundreds or
 219 thousands of features significantly contribute to a
 220 prediction, it is not reasonable to expect any user to
 221 comprehend why the prediction was made, even if
 222 individual weights can be inspected". Characteristics on
 223 which we generally believe it would be unethical to base
 224 such decisions—such as an individual's race or sex—
 225 may play a role in generating the ADM's decision
 226 without our knowledge. Even when such protected
 227 information is explicitly eliminated from the dataset,
 228 opaque AI technology may still display
 229 incomprehensible discrimination or 'prejudice by
 230 proxy.'² An ADM may discover a highly efficient method
 231 that utilizes a combination of factors (such as zip code
 232 and *alma mater*) in such a way that the output is
 233 tantamount to a judgment based on race. The more
 234 opaque the AI technology, the less certain we can be that
 235 it will be adequately unbiased in its assessment.

236 In response to the problems that can be generated by
 237 the use of opaque AI, there has been a general push for
 238 increasing transparency in AI. Governing bodies,
 239 technology watchdog groups, and ethicists have made
 240 transparency a priority in AI regulations. The European
 241 Commission's 2019 Ethics Guidelines for Trustworthy AI
 242 identifies transparency as its fourth out of seven key
 243 requirements that AI systems should meet. In January
 244 2020, the White House released its first guidelines for AI
 245 regulation which, although they are limited to the private
 246 sector and do not mention transparency verbatim, do
 247 include "trustworthiness," which is intimately connected
 248 to the value of transparency. Similarly, The Future of Life
 249 Institute explicitly includes two transparency-related
 250 items in their (2017) account of the general Principles of
 AI.³

245 Corporations such as Google and Microsoft have publicly
 246 acknowledged the importance of transparency in AI
 247 as

² See Barocas (2018).

³ These principles concern *failure transparency* (if an AI system causes harm, it should be possible to ascertain why), and *judicial*

transparency (any involvement by an autonomous system in judicial decision-making should provide a satisfactory explanation auditable by a competent human authority).

¹ <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scrap-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>.

well. As Microsoft CEO Satya Nadella stated in 2016, “We want not just intelligent machines but intelligible machines. Not artificial intelligence but symbiotic intelligence... People should have an understanding of how the technology sees and analyzes the world.” And in the framework for a ‘Good AI Society’, Floridi et al. (2018) call for enhanced explicability in AI when AI is involved in socially significant decisions. “Central to this framework is the ability for individuals to obtain a factual, direct, and clear explanation of the decision-making process, especially in the event of unwanted consequences” (p. 702). The consensus that seems to have emerged in response to the opacity problem has been to treat transparency in AI as valuable in and of itself, and that the overall benefit we gain from AI increases as transparency increases. That is, we are better off ethically the more transparent we make our AI models.

4 Ethically significant contexts: concerns beyond bias and fairness

Not all uses of opaque AI give rise to ethical concerns. There are many contexts in which the opacity of an AI model is insignificant simply because we consider the consequences of decisions made in those areas to be trivial. Intuitively, if certain activities genuinely qualify as “for entertainment purposes only,” such a context would be trivial, or at least not ethically significant. In the most general terms, for a context to be ethically significant the consequences of actions or decisions in that context must at minimum carry a risk of harm (where harm is very broadly construed).⁴

Robbins (2019) is skeptical of the call for transparency in AI, and suggests that while the use of opaque AI is ethically permissible in trivial contexts and certain non-trivial contexts (which he groups together as ‘neutral contexts’), it should not be allowed to operate in what he labels ‘morally sensitive contexts’.

Robbins intends this division between morally sensitive contexts and ‘neutral contexts’ to largely map onto the distinction between contexts in which we intuitively feel comfortable with the use of opaque AI and contexts in which this opacity seems potentially problematic. Commonly identified ethically problematic contexts of use are those such as judicial sentencing (Berk et al. 2016; Barry-Jester et al. 2015), predictive policing (Ahmed 2018; Ensign et al. 2017; Joh 2017; O’Neil 2016) and medical diagnosis (de Bruijne

2016; Dhar and Ranganathan 2015; Erickson et al. 2017). He writes,

One reason that using inexplicable decisions in morally sensitive contexts like the ones listed above is wrong is that we must ensure that the decisions are not based on inappropriate considerations... Combine this fact with using ML algorithms for decisions that have moral significance (i.e. decisions which could result in harm—broadly construed to include rights violations) and we have an ethically problematic situation. An algorithm used, for example, to accept or reject your loan request will significantly affect you. A rejection could cause you and your partner significant distress and change the course of your life. (Robbins, 2019, p. 498)

Robbins’s analysis seems to suggest that there are two features of a context which together make it a morally sensitive context. One concerns fairness. The other is magnitude of impact, or whether it is a “high-stakes” context. Regarding fairness, there is wide consensus that certain personal characteristics are ethically protected characteristics; these characteristics ought not be taken into account in high-stakes contexts—when the outcome of the decision can have a great impact on one’s welfare. Loan approval decisions, hiring decisions, recidivism risk and suitability for parole all seem to be areas in which we need to pay special attention to how judgments are made because there are fair and unfair procedures for making these judgments.

Given that there are clear cases in which we do and should value fairness over efficiency, and that it seems reasonable to interpret being treated unfairly as a kind of harm, contexts in which judgments might be made unfairly should be considered a type of high-stakes context with a significant risk of harm. If so, we can incorporate considerations of fair treatment in a general account of contexts in which there is significant risk of harm. Unfair treatment is one among many potential harms that we risk when we employ opaque AI; I propose that we widen the category of domains in which we might be prohibited from using opaque AI beyond those which fit Robbins’s description of “morally sensitive contexts” to include any context in which there is an opportunity to substantially impact the welfare or wellbeing of an individual or group. We can call these “ethically significant” contexts of use. Insofar as actions or decisions made in these areas can have significant impact on our wellbeing, special attention ought to be paid to our methods for arriving at decisions and determining our course of action in these areas. We may be ethically prohibited, for instance, from using opaque AI in hiring decisions because that model may exhibit unfair gender or racial bias, which has a significant impact on the welfare of those applicants. In the same way, we may be prohibited from using opaque

⁴Harm here is broadly construed to include (at minimum) opportunity costs, as well as intangible/unquantifiable harms such as rights violations, insufficient or inaccurate representation, harm to social reputation, and harm to self-esteem.

AI technology when deciding on actions regarding global food production: because the stability and resilience of the global food chain has a significant impact on human welfare, we may be ethically required to ensure that we have adequate understanding of the tools and processes on which we base those decisions.

The boundaries for what qualifies as an ethically significant context on my account are wide and somewhat vague, and may cast a wider-than-expected net over contexts that qualify as “ethically significant.” I believe the vagueness and breadth of this category accurately reflect the fact that our actual judgements regarding what features of the world qualify as ethically significant are notoriously difficult to codify.⁵ While these judgments are sometimes unpredictable, there are also central cases on which all or nearly all can agree. Additionally, unlike Robbins, I am not suggesting a blanket prohibition against the use of opaque AI in all ethically significant contexts. Therefore, identifying which specific cases qualify as ethically significant will not ultimately determine whether it is ethical to employ opaque AI in such a case. Rather, identifying a context as ethically significant means that we are required to subject that case to further scrutiny before we can determine whether it is ethical to employ opaque AI.

As indicated above, a more complete account of the costs and benefits of prohibiting the use of opaque AI in certain contexts will consider contexts beyond those in which issues of bias may arise. A more inclusive (but still incomplete) account of ethically significant contexts will include contexts in which there are multiple types of opportunity cost: risk of inappropriately skewed distribution of benefits (increasing inequity) as well as risk of missed opportunity for significant benefit (especially for vulnerable populations). Recognizing these features as relevant to the ethical significance of a situation allows us to treat cases in which opaque AI may be utilized in areas such as climate science, extreme weather event prediction, public health and medicine, [footnote here which says: see London (2019); Vincent (2018)] and global food production as ethically significant contexts. These areas have sometimes been misidentified as areas in which ethical concerns regarding AI opacity do not arise, because it seems obvious that we value efficiency over transparency in such cases.⁶ However, granting that we do in fact value efficiency over transparency in these areas does not entail that we cease to value transparency here, and it certainly does not entail that decisions and actions in these areas are ethically neutral or trivial. It would be a mistake to regard

areas in which our concern for efficiency wins out over our concern for transparency as areas that are “ethically neutral”, as Robbins (2019) seems to do. There are certain domains in which we value efficiency over transparency for *ethical* reasons, and to ignore this would grossly mischaracterize the domain of ethical concern. Rather, the particular ethical concerns we have in such cases are not put in sufficient jeopardy by the opacity of AI to justify the missed opportunity to substantially increase human welfare, which is itself a central ethical concern.

5 Transparency as a proxy for trustworthiness (or, If I knew what you know, I would not need to trust you)

An essential step toward answering the question of when, why, and to what extent we value transparency in AI is to identify the goal of increasing transparency. We can then ask whether that goal could be achieved by means other than transparency itself. Many have suggested that one of the main ethical goals⁷ in increased AI transparency is related to trust: we value transparency because it serves as a proxy for the trustworthiness of the AI model.

This is similar—but in at least one sense, importantly different—to the claim that, as transparency increases, stakeholders’ trust may reasonably increase as well.

Consider the domain of medical diagnostics. There is a widely supported movement for increased transparency in the AI tools that are currently used in making medical diagnoses, and the motivation behind the movement seems to be grounded in the importance of trust within the medical setting and the doctor-patient relationship. Trust and trustworthiness are two distinct but related features of that relationship, and both are essential to a successful expert/layperson relation. Whether a tool is trustworthy depends on the typical functioning of the tool—the actual overall predictive accuracy and reliability of the AI diagnostic tool, whether it is sufficiently robust in the face of small changes, and whether its predictions are based on a sufficiently broad and representative dataset. Trust, on the other hand, is a relation that holds between doctors and their diagnostic tools, or between doctors and the patients who rely on them. The presence of trust between doctor and patient increases the likelihood that the doctor will be able to effectively treat the patient; ideally, this improves the patient’s health-related wellbeing. This trust is appropriate—when it is—in part because society has guidelines in place to ensure

⁵ See Skerker, Purves, and Jenkins (2015) on the anti-codifiability

problem in robot and machine ethics.

⁶ See Robbins (2019) on valuing efficiency *rather than* transparency in certain non-trivial cases.

⁷ There are epistemic advantages to increasing transparency in AI models, but for the sake of this paper we are focusing solely on the ethical goals of requiring transparency in AI.

438 that a doctor's extensive training results in sound medical
439 judgment, and a well-functioning social system for
440 verifying expertise (such as board certification and
441 licensing). Trust is an essential feature of modern
442 society's successful (when it is successful) division of
443 epistemic labor. It is clearly indis- pensable for a
444 successful doctor-patient relationship, and the same
445 holds for the epistemic and ethical relationships between
446 experts and laypeople in general. Trust is essential in the
447 absence of understanding and explanation (with suf-
448 ficient understanding and explanation, trust can be
449 unnec- essary). It is often thought that we trust processes
448 that we understand, as Riberio et al. (2016) make explicit
449 here:

450 Whether humans are directly using machine
451 learning classifiers as tools, or are deploying
452 models within other products, a vital concern
453 remains: if the users do not trust a model or a
454 prediction, they will not use it. It is important to
455 differentiate between two different (but related)
456 definitions of trust: (1) trusting a predic- tion, i.e.
457 whether a user trusts an individual prediction
458 sufficiently to take some action based on it, and (2)
459 trusting a model, i.e. whether the user trusts a
460 model to behave in reasonable ways if deployed.
461 Both are directly impacted by how much the human
462 understands a model's behaviour, as opposed to
463 seeing it as a black box. (Riberio, 2016, section 1,
464 emphasis mine)

462 This is a common assumption regarding the relation
463 between trust and understanding, but it ignores the
464 addi- tional function and value of trust and
465 trustworthiness men- tioned above. Either increased trust
466 or increased understand- ing will typically result in an
467 agent's increased willingness to believe a certain
468 decision is accurate or engage with a certain tool. When
469 patients trust their doctors, that trust is not grounded in
470 the patients' understanding of the doctors' evidence or
471 reasoning. This remains opaque. Patients trust their
472 doctors because they know that, in a well-functioning
473 social system which includes institutions dedicated to
474 expert verification, a person would not hold the position
475 of doctor unless they possessed the adequate expertise.

474 In a society that operates with a successful division of
475 epistemic labor, trust and trustworthiness can replace
476 under- standing as epistemically and ethically sound
477 grounds for belief. Laypeople believe the judgments of
478 specialized experts because they trust those experts—
479 not because they understand their reasoning—and they
480 trust those experts because their social framework
481 includes institutions whose role it is to verify the
482 legitimacy of specialized experts. If the ultimate aim of
483 increased transparency is to establish trustworthiness
484 and build trust where appropriate, there may be other
485 avenues available for pursuing these goals—paths that
486 allow us to benefit from the power of opaque AI tech-
487 nology by verifying its trustworthiness. Transparency
488 itself need not be our goal.

486
487
488
489

If this is correct, then the options before us are either

(1) accept that the ethical concerns which give us reason
to employ opaque AI may outweigh the benefits of
transpar- ency, and determine how to best utilize opaque
AI given these epistemic limitations, or (2) refuse to
employ opaque AI in any ethically significant context on
the grounds that the use of an opaque process is ethically
impermissible in those contexts.

Given that there are enormous potential benefits
that could arise from the proper use of opaque AI in at
least some of the commonly identified ethically
significant domains— healthcare, climate science, the
global food chain, public safety—we would need
powerful ethical reasons to sup- port fully eliminating
its use in these areas. The success of the expert/layperson
division of epistemic labor shows us that many of our
ordinary, ethically responsible, and reli- able social
practices already implicitly reject (2) above: we routinely
employ opaque processes in ethically significant
domains. And I will argue that there is no special reason
to embrace (2) in the case of AI while rejecting it in the
case of human experts. If this is correct, then we are left
with option (1), and the ethical question before us is no
longer whether we ought to allow opaque AI to operate in
any ethically sig- nificant domains but rather what are
the most ethical ways of harnessing opaque AI in these
domains.

1 The expert/layperson relation—a blueprint for ethical opaque AI

I have suggested that we take our successful social
practice of deferring to specialized experts as a guide for
develop- ing an epistemically and ethically sound
method for utiliz- ing opaque AI. To this end, we will
need to examine when (i.e., under what conditions) it is
epistemically and ethically responsible to defer to experts
rather than relying on one's own reasoning. We also need
to know what features make an individual a genuine
expert, how, as a society, we determine that some
individual is an expert, and what methods we use for
deciding how to act when multiple experts disagree in
their decisions. Fortunately, these questions are
beginning to receive increased attention both in sociology
and philosophy, under the general headings of social
epistemology and the epistemology of testimony.
Goldman 2001; Goldman 2014; Lackey 2016;

In what follows I will make a preliminary case for the
claim that the essential features of experts-- the features
that make expert opinion trustworthy, and our trust in
those individuals' decisions both epistemically and
ethically responsible—can be realized in AI as well. For
this to be the case, the relevant features of human experts
must not be essentially human features. Certainly,
human experts have noteworthy features that AI models
lack; for instance, we

typically assume that human experts have a concept of the greater good and a desire to promote it. If such traits play an indispensable role in generating the trust and trustworthiness on which the expert/layperson relation depends, this relation will not be a viable model for the ethical use of opaque AI. As I hope to show below, the trust that exists in the expert/layperson relation is not fundamentally based on faith in the moral goodness of the expert but rather on the nature of expertise and the existence of institutions that serve to verify these experts. If these features are not uniquely human features, then, insofar as we have ethically acceptable methods of evaluating when we ought to defer to human experts in high-stakes contexts, we have a potential framework for determining when it is ethically appropriate to defer to the decisions generated by opaque AI technology.

In the mid 1980's, philosopher John Hardwig sparked renewed interest in the social aspects of knowledge-building by drawing attention to the myriad situations in which we are better off—rationally speaking—deferring to someone else's judgment on a particular matter rather than attempting to reason through that matter ourselves. These are situations in which the matter at hand concerns an area of highly specialized knowledge, and there are highly trained experts who specialize in that area. In such a case, a layperson would be more rationally justified in deferring to the expert's judgment than they would in performing their own independent reasoning and standing by the judgment at which they themselves had arrived. That is to say, a layperson has better reasons to believe an expert's judgment is correct than his or her own, even when that judgment conflicts with theirs. Assuming that the layperson is a genuine layperson, and the expert a genuine expert, Hardwig writes,

If, then, layman B (1) has not performed the inquiry that would provide the evidence for his belief that p, (2) is not competent, and perhaps could not even become competent, to perform that inquiry, (3) is not able to assess the merits of the evidence provided by expert A's inquiry, and (4) may not even be able to understand the evidence and how it supports A's [the expert's] belief that p, can B nonetheless have good reasons to believe that A has good reasons to believe that p? I think he can. If so, should we conclude that B's belief that p is rationally justified? I think we should, acknowledging that B's belief stands on better epistemic ground than other beliefs which we would call simply irrational or nonrational. (1985, p.339)

Following Hardwig, we can say that in order for laypeople to be justified in deferring to the (opaque) reasoning of experts—rather than being rationally required to perform their own (transparent) reasoning—there are (at least) three criteria that must be met[R1].

1. The laypeople have not, themselves, performed the reasoning that is being left to the expert.
2. The laypeople are not capable of performing the reasoning that is being left to the expert (for any of several possible reasons, to be discussed below).
3. The laypeople cannot themselves 'assess the merits of the evidence' nor understand how the evidence supports the expert's decision. (This combines 3 and 4 in Hardwig's criteria, above).

6.1 Ruling in-and ruling out-the use of opaque AI

As will soon become apparent, even a framework intended to show where we are permitted to employ opaque AI technology in ethically significant contexts will rule *against* the use of opaque AI in many of the notoriously problematic cases in which those models are already in use. Below, I will adapt Hardwig's (minimal) criteria for deference to experts to apply to AI and briefly discuss the most readily apparent implications of interpreting each criterion in these particular ways.

1. Neither transparent models nor humans have performed the task in question on the scale at which the opaque AI model will be performing that task.

Our general motivation for applying AI to any particular task becomes more clear when we draw attention to the scale of the task; additionally, explicitly specifying the scale of the task is essential to properly characterizing the task itself. In broad terms, many of the same types of tasks that AI models are designed to perform—reviewing loan applications, evaluating job candidates, deciding how to deploy police resources, predicting effects of climate and weather events on food production—have all previously been performed by human reasoners. But the size of the problems to which we might apply the tools of AI, and the scale on which we intend for these tasks to now be performed, is unprecedented. These tasks may require more labor-hours than we can reasonably expect from human beings, especially when the tasks are time-sensitive.

That said, if this first criterion must be met for any ethically responsible application of opaque AI in an ethically significant context, then many instances in which opaque AI is already being deployed may not satisfy the criteria necessary for the ethical use of opaque AI. (More will be said about this when we discuss guideline (B) in the following section.)

2. Transparent models are not practically capable of performing the task that the opaque AI model is intended to perform.

Whether this criterion is met will in part depend on the state of AI technology and the actual skillsets of AI researchers at the time the decision is being made. Rudin (2019) points to this aspect of the problem when she writes,

Black box models seem to uncover ‘hidden patterns’. The fact that many scientists have difficulty constructing interpretable models may be fueling the belief that black boxes have the ability to uncover subtle hidden patterns in the data about which the user was not previously aware. A transparent model may be able to uncover these same patterns. If the pattern in the data was important enough that a black box model could leverage it to obtain better predictions, an interpretable model might also locate the same pattern and use it. Again, this depends on the ML researcher’s ability to create accurate yet interpretable models. The researcher needs to create a model that has the capability of uncovering the types of pattern that the user would find interpretable, but also the model needs to be flexible enough to fit the data accurately. This, and the optimization challenges discussed above, are where the difficulty lies with constructing interpretable models. (2019, p. 201)

If equally proficient transparent models⁸ already exist or could realistically be developed within the requisite timeframe (where ‘equally proficient’ takes into account the *speed* required to perform the task effectively as well as the scale of the task), the additional value conferred by their transparency may make them ethically preferable to an opaque model. Though Rudin is optimistic regarding the potential of transparent (in this case, interpretable) models to perform as well as opaque models, this is by no means guaranteed. As she acknowledges, “This problem is compounded by the fact that researchers are now trained in deep learning, but not in interpretable ML...” and “It could be possible that there are application domains where a complete black box is required for a high stakes decision,” though she notes that, “As of yet, I have not encountered such an application” (2019, p. 207).

⁸ While “opaque” has a standard meaning in the literature on this topic, “transparent” has several common meanings when used in the context of AI models. A satisfactorily transparent AI model might be an interpretable model, or an explicable model, or it may be comprehensible to the relevant practitioner or stakeholder, etc. A thorough account of how “transparency” has been interpreted in the literature on AI regulations is beyond the scope of this discussion, but see Chen 2018; Li et al. 2018; Lipton 2016; Miller 2017; Mittelstadt et al. 2019; Molnar 2019; Riberio 2016; Rudin 2019; Zerilli 2002;

3. e are unable to satisfactorily explain the AI model within a reasonable amount of time given the urgency of the task in question.

An explanation of an AI model would allow us to “assess the merits” of the evidence on which the model

is basing its decision and “understand... how [the evidence] supports” that decision. The third criterion roughly specifies that in order for us to sacrifice transparency for the benefits gained by employing opaque AI in a particular ethically significant context, that opacity must be a result of our genuine inability to explain the operations of the AI model, rather than an unwillingness to deploy sufficient resources to the task. (Note that this issue will only arise when there is a question of irresponsibly employing opaque AI—the context itself must be ethically significant for ethical concerns to compete with the value of the efficiency or accuracy gained by utilizing opaque AI models.)

In addition to this cursory description of when it would be reasonable for a layperson to defer to the judgment of an expert, Hardwig also provides a rough approximation of the personal features that make an individual an expert. Briefly, an expert must have engaged in “inquiry that has been sustained, prolonged, and systematic” (1985, p. 338). Though we would need to determine what features of an AI model would make its “inquiry” into a specific domain suitably “sustained, prolonged, and systematic,” this criterion seems to pose no special difficulty for AI. And given that these models fundamentally function by discovering and attuning themselves to patterns in data, such data-processing operations should satisfy all relevant features of an “inquiry.”

6.2 The social institutions/practices underwriting our successful deference to experts (and how they might be replicated in the case of AI)

So far I have proposed a preliminary set of criteria that would need to be met in order for an individual—or an AI model—to qualify as an expert, as well as conditions under which may it be epistemically and ethically responsible to defer to the judgments of a human or artificial “expert”. In this section, we will consider preliminary ideas regarding how we might determine whether some opaque AI model should be considered an expert in this sense. An opaque model may possess the requisite features for “expertise” in a certain area, but the opacity of that model will make it challenging for us to know whether the model has satisfied the appropriate criteria. In addition, I will make preliminary suggestions for how we might deal with morally weighty cases in which (just as with human experts) multiple opaque AI models disagree in their predictions or decisions.

In the familiar cases of human experts, the answers to both of these questions rely, in part, on the existence of

727 a larger network of experts in addition to the individual
 728 (potential) expert in question. In areas of technical
 729 speciali- zation, (academic research, professions such as
 730 journalism and law, etc.) we commonly find established
 731 institutions and professional organizations that grant
 732 degrees, credentials, or otherwise certify that the
 733 individual in question does in fact qualify as an expert.
 734 These organizations are typically com- posed of
 735 individuals who themselves possess certain types of
 736 relevant expertise. When cases arise in which a purported
 737 ‘expert’ fails to meet the standards set by the certifying
 738 bod- ies in their fields, we rely on these institutions to
 739 revoke that individual’s credentials. Lawyers can be
 740 disbarred, doctors can lose their license to practice
 741 medicine, journalists can lose their press credentials, and
 742 so on. Ideally, this process serves to inform the public that
 743 these individuals are not, in fact, genuine experts in their
 744 supposed fields. These institu- tions allow laypeople to
 745 know which individuals are experts in which fields, and
 746 responsibly defer to their judgments, even though
 747 exactly what makes that individual an expert in that field
 748 is beyond the understanding of the layperson.

745 The presence of multiple experts within a single field
 746 is not only essential to our ability to know which
 747 individuals are experts (since we, as laypeople, cannot
 748 evaluate their expertise for ourselves); the fact that large
 749 numbers of inde- pendent experts regularly converge in
 750 their opinions give us an imperfect but reliable indication
 751 that these judgments are correct, as well as a means of
 752 determining how to act when experts disagree. If a
 753 significant majority of genuine experts converge in their
 754 opinion on a particular issue, and a small number of
 755 experts disagree, it will be reasonable for the layperson
 756 to accept the opinion of the majority.

755 Adapting our methods for certifying experts and
 756 handling expert disagreements such that we can apply
 757 them to opaque AI presents more of a challenge than
 758 adapting the criteria for expertise itself or for responsibly
 759 deferring to experts. The relationship between laypeople
 760 and experts in modern soci- ety has a long history, and
 761 the trustworthiness of these cre- dentialing institutions is
 762 born out only by society’s repeated knowledge-building
 763 success over time. Our engagement with opaque AI
 764 technology has both a short and checkered past. We do
 765 not have the convenience of a lengthy history—on a
 766 human timescale—to indicate which methods for
 767 certifying the expert-status of an opaque AI model will
 768 prove to be trustworthy, and which methods are likely to
 769 fail.

767 Because of the importance that time plays in revealing
 768 the reliability of expert decisions, of our method for ver-
 769 ifying individuals as genuine experts, and of our division
 770 of epistemic labor in general, whatever way in which we
 771 choose to adapt this feature to create an analogous
 772 method for revealing the trustworthiness of any opaque
 773 AI technol- ogy will be highly speculative. There are no
 774 obvious candi- dates for artificial analogs of the passage
 775 of time. With that in mind, one possible option would
 776 be to treat the notion

780 of an epoch in artificial neural networks as a stand-in for
 781 the ordinary passage of time. Rather than thinking of
 782 the history of AI models on a human timescale, it may be
 783 more appropriate to frame the notion of “an adequate
 784 length of time” on which to judge the reliability of an AI
 785 model to reflect an AI timescale. So whereas ANNs and
 786 other deep learning models may have emerged 10 years
 787 ago on a human timescale, a massive number of epochs
 788 for those models has passed within this span.
 789 (Determining the optimal number of epochs for training a
 790 neural network is currently considered something of an
 791 art in machine learning.)

790 There are a growing number of organizations
 791 dedicated to developing something akin to a
 792 “credentialing processes” for AI. The Institute of
 793 Electrical and Electronics Engineers (IEEE)
 794 continuously updates its standards for the devel-
 795 opment and use of AI. The International Organization
 796 for Standardization (ISO) and The International
 797 Electrotech- nical Commission (IEC) both work to
 798 develop standards that aim to make AI more “resilient,
 799 reliable, accurate, and secure”. And the European
 800 Commission’s 2021 proposal for Regulation on Artificial
 801 Intelligence includes a legal frame- work by which to
 802 judge the risk of AI. The UK Institute for Ethical AI and
 803 Machine Learning, the Global Partnership on AI (GPAI),
 804 and the OECD AI Policy Observatory all support projects
 805 and policy aimed at increasing trustworthiness in AI.
 806 What form a successful credentialing process will eventu-
 807 ally take, and to what extent these certification systems
 808 are already in place, is a question to be addressed
 809 elsewhere. But if we are interested in developing an
 810 approval process that could certify certain opaque AI
 811 technology and approve its use in particular contexts
 812 while allowing the technology to remain opaque, we
 813 might make progress on this issue by continuing
 814 research into the relevant features of familiar and
 815 successful practices of certifying human experts.

811 The final feature of the expert-layperson relationship
 812 that
 813 we will address here—our methods for dealing with cases
 814 of expert disagreement—is simple to adapt in theory
 815 (though perhaps less so in practice). Our successful
 816 division of epis- temic labor crucially depends on the
 817 existence of multiple independently trained experts in a
 818 single field, addressing the same issue and converging on
 819 the same opinion through a variety of independent
 820 methods. At the present moment, it is unclear whether
 821 there exists a sufficient number—and variety—of AI
 822 models that could perform the same ethically significant
 823 task (whatever this task may be) to deal with disa-
 824 greement in an analogous way. But there may be no better
 825 way to establish the requisite level of trustworthiness⁹ [1]
 826 of an opaque AI model than developing multiple,
 827 independent,

⁹ to whatever extent is required such that it would be ethically
 responsible to utilize that opaque technology in the particular ethi-
 cally significant context in question.

828 opaque models, operating with distinct architecture
829 and trained on distinct (but appropriately relevant) data
830 sets, and finding that they converge on the same decision.
831 Given that opaque AI will be an ever-present ethical
832 issue, developing multiple models to perform the same
833 ethically significant task may be well worth the
investment.

834 2 Preliminary guidelines for the ethical use 835 of opaque AI

836 Given that I claim it may be ethically permissible
837 (perhaps required) to use opaque AI technology in
838 certain ethically significant contexts, this section
839 provides a plausible decision procedure for evaluating
840 whether a particular context is one in which we could
841 ethically employ opaque AI. I suggest three general
842 questions that should be addressed in the process of
making such a decision.

- 843 (A) Is the context in question an ethically significant
844 con- text?
845
846 (B) Could the task at issue be performed equally well by
847 a transparent process?
848
849 (C) Are the benefits of successfully performing this task
850 greater than both (i) the cost of potentially failing at
851 this task (whatever constitutes “failure” in this case)
and (ii) the cost of not performing this task at all?

852 (A) Is the context in question an ethically significant
853 con- text? The process of evaluation begins with question
854 (A): Is the context in question an ethically significant
855 context? If we can be reasonably certain that the answer to
856 (A) is “no,” then the ethical concerns surrounding the
857 use of opaque AI do not arise in this situation, and we are
858 at liberty to use opaque AI for the task at issue. Note that
859 the triviality or ethical significance of a context will most
860 often be decided according to a broad and diverse set of
861 standards, some of which may involve apparently
862 objective and quantifiable measures (for example, the
863 potential consequences of utilizing some proposed AI
864 technology in the global food supply chain) and some of
865 which may involve standards that will vary relative to a
866 cultural or social context (the impact of utilizing some
867 proposed AI on the representation of a particular socially
868 marginalized or vulnerable group). Note also that the
869 ethical significance of a context will be a matter of degree,
870 depend- ing on the gravity of the particular situation(s)
871 involved. If the answer to (A) is “yes,” then we need to
872 address question (B). Could the task at issue be
873 performed equally well by a transparent process
874 (whether human or AI)? This question will be familiar
875 from the criteria for rationally deferring to experts in
general. The additional benefits that arise from
transparency in how decisions are made in all ethically sig-
nificant contexts may outweigh whatever benefits the
opaque

AI may provide. Here, it is important to note that “per-
forming a task equally as well” will include—at
minimum— issues of equity and fairness in addition to
efficiency and accuracy. As noted in Sect. 4 3, we cannot
entirely ignore the harms of opportunity costs for the
sake of eliminating bias, especially when those costs are
borne by marginalized and vulnerable populations. To
permit the use of opaque AI in an ethically significant
context, the answer to question

(A) must be yes, and the answer to question (B) must be
no. If so, then it may be ethically permissible to utilize
opaque AI, subject to further consideration, such as those
raised in question (C). Are the benefits of successfully
performing this task greater than both (i) the cost of
potentially failing at this task (whatever constitutes
“failure” in this case) and

(ii) the cost of not performing this task at all? If there are
ethically significant cases in which all three bars are met,
then there are non-trivial cases in which we would be
per- mitted—perhaps required—to utilize opaque AI.
And given that meeting all three bars requires that the
opaque model in question be reliable and trustworthy, we
will need a frame- work for evaluating the reliability and
trustworthiness of opaque AI technology. I hope to have
made a preliminary case for looking to our successful
social practice of deferring to experts in ethically
significant domains for a blueprint of how to responsibly
employ opaque AI in such a case.

3 Conclusion

I acknowledge that, even as guidelines go, those given
above are considerably vague. I view this vagueness as
appropri- ate, and—practically speaking—ineliminable.
Here, we are concerned with developing rules for ethical
action in the use of AI, and as Aristotle said, we should
only look for preci- sion in each class of things just so far
as the nature of the subject admits. Any rule, no matter
how precise, requires interpretation when applied to a
particular case. And when the interpretation of those
rules involves disentangling and weighing competing
moral values, it is the process of inter- pretation itself—
and not the rule—that will be doing the lion’s share of
the work. So I would suggest that insofar as these
guidelines are vague, their vagueness is appropri- ate to
the subject at hand. Deciding whether a task could be
performed equally well by some satisfactorily transpar-
ent (human or algorithmic) decision-making process
will involve weighing competing values, and the relative
strength of those competing values will depend on the
ethical inclina- tions of the individuals performing the
evaluation. There is no standard, universally applicable
measure for assigning weights to these values; each case
will need to be evalu- ated individually, and an argument
will need to be made for weighting any of these values
more strongly than the others. The same is true for
deciding whether the benefits of success

926 are worth the potential costs of failure. Human
927 judgment cannot be entirely removed from decision-
928 making in ethically significant domains, no matter how
929 trustworthy the AI technology involved. At minimum,
930 humans must still be in-the-loop to (1) make case-
931 specific value-judgments, and

932 (2) make cost/benefit assessments in cases where the
933 costs and benefits are not fully commensurable. And
934 given that we are discussing opaque AI technology,
935 humans will need to be in-the-loop to monitor for
936 potential instances of biased outcomes. The threat of bias
937 will remain, whether or not the cost of that potential bias
938 is outweighed by the potential benefits of a successful
939 outcome.

940 These guidelines are not intended to serve as a
941 complete checklist for the ethical use of opaque AI. They
942 merely offer one plausible set of rules for evaluating
943 whether some situation is an instance in which we
944 should consider, or refuse, to employ certain opaque AI
945 technology. If we decide we should, we might then look
946 to the blueprint provided by the expert/layperson
947 division of epistemic labor to see how to do so well. In
948 addition, the overview of the expert/layperson relation
949 given above is not intended to fully capture the robust
950 and complex features of this social epistemic practice.
951 Whether this overview accurately represents the
952 fundamental features of this relationship is separate
953 from the question of whether the expert/layperson
954 relation itself—and the institutions that support it—can
955 provide us with a general framework for developing an
956 ethical approach to harnessing the power of opaque AI,
957 as I believe it can.

958 **Funding** No funding was received to assist with the preparation of
959 this manuscript.

960 Declarations

961 **Conflict of interest** On behalf of all authors, the corresponding
962 author states that there is no conflict of interest.

963 References

- 964 Ahmed M (2018) Aided by Palantir, the LAPD uses predictive polic-
965 ing to monitor specific people and neighborhoods. *The*
966 *Intercept*. [https://theintercept.com/2018/05/11/predictive-](https://theintercept.com/2018/05/11/predictive-policing-surveillance-ce-los-angeles/)
967 [policing-surveillance-ce-los-angeles/](https://theintercept.com/2018/05/11/predictive-policing-surveillance-ce-los-angeles/)
968 Barocas S (2018) Accounting for artificial intelligence: rules,
969 reasons, rationales. In: Human rights, ethics, and artificial
970 intelligence, 30 Nov. Harvard Kennedy School Carr Center for
971 Human Rights Policy. Lecture
972 Barry-Jester A, Casselman B, Goldstein D (2015) The new science
973 of sentencing. *The Marshall Project*.
974 [https://www.themarshallproject.](https://www.themarshallproject.org/2015/08/04/the-new-science-of-sentencing)
975 [org/2015/08/04/the-new-](https://www.themarshallproject.org/2015/08/04/the-new-science-of-sentencing)
976 [science-of-sentencing](https://www.themarshallproject.org/2015/08/04/the-new-science-of-sentencing)
977 Berk RA, Sorenson SB, Barnes G (2016) Forecasting domestic vio-
978 lence: a machine learning approach to help inform
979 arraignment decisions. *J Empir Leg Stud* 13(1):94–115.
980 <https://doi.org/10.1111/jels.12098>

- 981 Chen C et al (2018) This looks like that: deep learning for interpret-
982 able image recognition. Preprint at [https://arxiv.org/abs/1806.](https://arxiv.org/abs/1806.10574)
983 [10574](https://arxiv.org/abs/1806.10574)
984 de Bruijne M (2016) Machine learning approaches in medical
985 image analysis: from detection to diagnosis. *Med Image Anal*
986 33:94–97. <https://doi.org/10.1016/j.media.2016.06.032>
987 Dhar J, Ranganathan A (2015) Machine learning capabilities in medi-
988 cal diagnosis applications: computational results for
989 hepatitis dis- ease. *Int J Biomed Eng Technol* 17(4):330–340.
990 <https://doi.org/10.1504/IJBET.2015.069398>
991 Erickson BJ, Korfiatis P, Akkus Z, Kline TL (2017) Machine
992 learning for medical imaging. *Radiographics* 37(2):505–515.
993 <https://doi.org/10.1148/rg.2017160130>
994 Ensign, D., Friedler, S. A., Neville, S., Scheidegger, C., & Venkatasu-
995 bramanian, S. (2017). Run away feedback loops in predictive
996 policing. In *Proceedings of machine learning research*, 81, 1–
997 12. Retrieved from <http://arxiv.org/abs/1706.09847>
998 European Commission (2019) Ethics Guidelines for Trustworthy
999 AI. <https://ec.europa.eu/futurium/en/ai-allianceconsultation>.
1000 Floridi L, Cowls J, Beltrametti M, Chatila R, Chazerand P, Dignum
1001 V, Vayena E (2018) AI4People—an ethical framework for a
1002 good AI society: opportunities, risks, principles, and recom-
1003 mendations. *Mind Mach* 28(4):689–707.
1004 <https://doi.org/10.1007/s11023-018-9482-5>
1005 Future of Life Institute (2017) Asilomar AI Principles. <https://futureoflife.org/ai-principles/>
1006 Goldman AI (2001) Experts: which ones should you trust? *Philos*
1007 *Phenomenol Res* 63(1):85–110
1008 Goldman AI (2014) Social process reliabilism: solving justification
1009 problems in collective epistemology. *Lackey* 2014:11–41.
1010 <https://doi.org/10.1093/acprof:oso/9780199665792.003.0002>
1011 Günther M, Kasirzadeh A (2022) Algorithmic and human decision
1012 making: for a double standard of transparency. *AI & Soc.*
1013 <https://doi.org/10.1007/s00146-021-01200-5>
1014 Hardwig J (1985) Epistemic dependence. *J Philos* 82(7):335–349
1015 Joh, E. E. (2017). Feeding the machine: Policing, crime data, & algo-
1016 rithms. *William & Mary Bill of Rights Journal*, 26, 287.
1017 Lackey J (2016) What is justified group belief? *Philos Rev* 125(3):341–
1018 396. <https://doi.org/10.1215/00318108-3516946>
1019 Li O, Liu H, Chen C, Rudin C (2018) Deep learning for case-based
1020 reasoning through prototypes: a neural network that explains
1021 its predictions. In: *Proceedings of AAAI Conference on*
1022 *Artificial Intelligence* 3530–3537 (AAAI, 2018).
1023 Lipton ZC (2016) The mythos of model interpretability. In: *ICML*
1024 *Workshop on Human Interpretability in Machine Learning*,
1025 vol 2017, pp. 96–100, 24
1026 London AJ (2019) Artificial intelligence and black-box medical
1027 deci- sions: accuracy versus explainability. *Hastings Center*
1028 *Rep* 49. <https://doi.org/10.1002/hast.973>
1029 Miller T (2017) Explanation in artificial intelligence: insights from
1030 the social sciences. *arXiv*
1031 Mittelstadt B, Russell C, Wachter S (2019) Explaining explanations
1032 in AI. In: *Proceedings of fairness, accountability, and*
1033 *transparency (FAT*) (ACM, 2019)*
1034 Molnar C (2019) Interpretable machine learning
1035 Nadella S (2016) Microsoft's CEO explores how humans and AI
1036 Can solve society's challenges—together. *Slate*.
1037 [https://slate.com/techn](https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html)
1038 [ology/2016/06/microsoft-ceo-satya-](https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html)
1039 [nadella-humans-and-a-i-can-](https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html)
1040 [work-together-to-solve-](https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html)
1041 [societyschallenges.html](https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html)
1042 O'Neil C (2016) Weapons of math destruction: how big data
1043 increases inequality and threatens democracy. *Crown*
1044 Ribeiro MT, Singh S, Guestrin C (2016) “Why should I trust you?":
1045 explaining the predictions of any classifier. *KDD*
1046 Robbins S (2019) A misdirected principle with a catch: explicabil-
1047 ity for AI. *Mind Mach* 29:495–514. [https://doi.org/10.1007/](https://doi.org/10.1007/s11023-019-09509-3)
1048 [s11023-019-09509-3](https://doi.org/10.1007/s11023-019-09509-3)

- 1042 Rudin C (2019) Stop explaining black box machine learning models
1043 for high stakes decisions and use interpretable models instead.
1044 Nat Mach Intell 1:206–215
- 1045 Skerker M, Purves D, Jenkins R (2015) Autonomous machines,
1046 moral judgment, and acting for the right reasons. Ethical
1047 Theory Moral Pract 18(4):851–872 (**Special Issue: BSET-
1048 2014**)
- 1049 Vincent J (2018) AI that detects cardiac arrests during emer- gency
1050 calls will be tested across Europe this summer. The Verge.
1051 [https:// www. theverge. com/ 2018/4/ 25/ 17278 994/ ai-
1052 cardiac-arrest-corti-emergency-call-response](https://www.theverge.com/2018/4/25/17278994/ai-cardiac-arrest-corti-emergency-call-response)
- 1053 Zerilli J (2022) Explaining machine learning decisions. Philos Sci
1054 89(1):1–19
- 1055 Zerilli J, Knott A, Maclaurin J et al (2019) Transparency in
1056 algorithmic and human decision-making: is there a double
1057 standard? Philos Technol 32:661–683.
1058 <https://doi.org/10.1007/s13347-018-0330-6>
- Publisher's Note** Springer Nature remains neutral with regard to
jurisdictional claims in published maps and institutional
affiliations.
- Springer Nature or its licensor holds exclusive rights to this article
under a publishing agreement with the author(s) or other
rightsholder(s); author self-archiving of the accepted manuscript
version of this article is solely governed by the terms of such
publishing agreement and applicable law.

REVISED PROOF