



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

Amber Ross<sup>1</sup>

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

✉ Amber Ross  
amber.ross@ufl.edu

<sup>1</sup> Department of Philosophy, The University of Florida,  
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 in  
 162 specialized fields (making the required reasoning opaque to  
 163 most), AI is most useful in areas where its speed and capac-  
 164 ity for data processing greatly surpasses human abilities—  
 165 the same factors that make certain AI technology opaque.  
 166 And just as the judgment of experts is most valuable in  
 167 high-stakes situations, the maximal benefit we can derive  
 168 from AI will be in its application to areas that are central to  
 169 human welfare (areas such as health, agriculture, climate,  
 170 and public safety). The power of AI is a double-edged sword.  
 171 Its extraordinary speed and unconventional data-processing  
 172 methods are the same factors that can make the most power-  
 173 ful AI opaque to its users and stakeholders, creating ethical  
 174 concerns regarding whether it ought to be used in the very  
 175 areas in which it could potentially provide the most benef-  
 176 it. The more knowledge we are ethically required to have  
 177 regarding how AI technology works when it operates in a  
 178 particular domain, the less likely it is that we will be ethi-  
 179 cally permitted to use AI in that domain.

### 180 3 Opacity and the problem of algorithmic 181 bias

182 The call for transparency in AI aims at safeguarding and  
 183 improving human welfare—in particular, by protecting  
 184 vulnerable groups who are most often harmed by opaque  
 185 AI technology and marginalized in AI development. This  
 186 goal is and should be a top priority in AI regulation. The  
 187 speed and processing power of AI not only comes at an epis-  
 188 temic cost; as we have learned, our limited epistemic access  
 189 to certain AI models can bring with it ethical costs as well.  
 190 In 2016, investigative journalists at ProPublica published  
 191 an article that exposed apparent racial bias in the popular  
 192 risk-assessment software COMPAS, used to aid judicial  
 193 decision-making regarding individuals' risk of recidivism  
 194 and eligibility for parole. In 2018, Reuters<sup>1</sup> revealed that  
 195 the AI hiring algorithm in development at Google showed  
 196 a strong gender bias.

197 The push to integrate these ADMs into areas such as  
 198 recidivism risk assessment, loan approval, and hiring prac-  
 199 tices, has exposed a tension between two worthwhile goals:  
 200 (i) increased efficiency in important decision-making pro-  
 201 cesses and (ii) protecting individuals' rights by ensuring  
 202 such decisions are based only upon ethically appropriate  
 203 considerations. This tension can become more problemat-  
 204 ic 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 or  
 207 stakeholder.

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

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

249 Corporations such as Google and Microsoft have  
 250 publicly acknowledged the importance of transparency  
 251 in AI as well. As Microsoft CEO Satya Nadella stated in

<sup>2</sup> See Barocas (2018).

<sup>3</sup> 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).

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

252 2016, “We want not just intelligent machines but intelligible  
 253 machines. Not artificial intelligence but symbiotic  
 254 intelligence... People should have an understanding of how  
 255 the technology sees and analyzes the world.” And in the  
 256 framework for a ‘Good AI Society’, Floridi et al. (2018)  
 257 call for enhanced explicability in AI when AI is involved in  
 258 socially significant decisions. “Central to this framework  
 259 is the ability for individuals to obtain a factual, direct, and  
 260 clear explanation of the decision-making process, especially  
 261 in the event of unwanted consequences” (p. 702). The  
 262 consensus that seems to have emerged in response to the  
 263 opacity problem has been to treat transparency in AI as  
 264 valuable in and of itself, and that the overall benefit we  
 265 gain from AI increases as transparency increases. That is,  
 266 we are better off ethically the more transparent we make  
 267 our AI models.

#### 268 4 Ethically significant contexts: concerns 269 beyond bias and fairness

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

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

286 Robbins intends this division between morally sensitive  
 287 contexts and ‘neutral contexts’ to largely map onto the dis-  
 288 tinction between contexts in which we intuitively feel com-  
 289 fortable with the use of opaque AI and contexts in which  
 290 this opacity seems potentially problematic. Commonly iden-  
 291 tified ethically problematic contexts of use are those such  
 292 as judicial sentencing (Berk et al. 2016; Barry-Jester et al.  
 293 2015), predictive policing (Ahmed 2018; Ensign et al. 2017;  
 294 Joh 2017; O’Neil 2016) and medical diagnosis (de Bruijne  
 295 2016; Dhar and Ranganathan 2015; Erickson et al. 2017).  
 296 He writes,

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

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

323 Given that there are clear cases in which we do and  
 324 should value fairness over efficiency, and that it seems  
 325 reasonable to interpret being treated unfairly as a kind of  
 326 harm, contexts in which judgments might be made unfairly  
 327 should be considered a type of high-stakes context with a  
 328 significant risk of harm. If so, we can incorporate consid-  
 329 erations of fair treatment in a general account of contexts  
 330 in which there is significant risk of harm. Unfair treatment  
 331 is one among many potential harms that we risk when we  
 332 employ opaque AI; I propose that we widen the category  
 333 of domains in which we might be prohibited from using  
 334 opaque AI beyond those which fit Robbins’s description of  
 335 “morally sensitive contexts” to include any context in which  
 336 there is an opportunity to substantially impact the welfare  
 337 or wellbeing of an individual or group. We can call these  
 338 “ethically significant” contexts of use. Insofar as actions or  
 339 decisions made in these areas can have significant impact  
 340 on our wellbeing, special attention ought to be paid to our  
 341 methods for arriving at decisions and determining our course  
 342 of action in these areas. We may be ethically prohibited, for  
 343 instance, from using opaque AI in hiring decisions because  
 344 that model may exhibit unfair gender or racial bias, which  
 345 has a significant impact on the welfare of those applicants.  
 346 In the same way, we may be prohibited from using opaque  
 347 AI technology when deciding on actions regarding global  
 348 food production: because the stability and resilience of the

<sup>4</sup> FL01 Harm here is broadly construed to include (at minimum) opportu-  
 FL02 nity costs, as well as intangible/unquantifiable harms such as rights  
 FL03 violations, insufficient or inaccurate representation, harm to social  
 FL04 reputation, and harm to self-esteem.

349 global food chain has a significant impact on human welfare,  
350 we may be ethically required to ensure that we have adequate  
351 understanding of the tools and processes on which we base  
352 those decisions.

353 The boundaries for what qualifies as an ethically signifi-  
354 cant context on my account are wide and somewhat vague,  
355 and may cast a wider-than-expected net over contexts that  
356 qualify as “ethically significant.” I believe the vagueness  
357 and breadth of this category accurately reflect the fact that  
358 our actual judgements regarding what features of the world  
359 qualify as ethically significant are notoriously difficult to  
360 codify.<sup>5</sup> While these judgments are sometimes unpredict-  
361 able, there are also central cases on which all or nearly all  
362 can agree. Additionally, unlike Robbins, I am not suggest-  
363 ing a blanket prohibition against the use of opaque AI in all  
364 ethically significant contexts. Therefore, identifying which  
365 specific cases qualify as ethically significant will not ulti-  
366 mately determine whether it is ethical to employ opaque  
367 AI in such a case. Rather, identifying a context as ethically  
368 significant means that we are required to subject that case to  
369 further scrutiny before we can determine whether it is ethical  
370 to employ opaque AI.

371 As indicated above, a more complete account of the costs  
372 and benefits of prohibiting the use of opaque AI in certain  
373 contexts will consider contexts beyond those in which issues  
374 of bias may arise. A more inclusive (but still incomplete)  
375 account of ethically significant contexts will include contexts  
376 in which there are multiple types of opportunity cost: risk of  
377 inappropriately skewed distribution of benefits (increasing  
378 inequity) as well as risk of missed opportunity for significant  
379 benefit (especially for vulnerable populations). Recogniz-  
380 ing these features as relevant to the ethical significance of  
381 a situation allows us to treat cases in which opaque AI may  
382 be utilized in areas such as climate science, extreme weather  
383 event prediction, public health and medicine,<sup>6</sup> and global  
384 food production as ethically significant contexts. These areas  
385 have sometimes been misidentified as areas in which ethi-  
386 cal concerns regarding AI opacity do not arise, because it  
387 seems obvious that we value efficiency over transparency  
388 in such cases.<sup>7</sup> However, granting that we do in fact value  
389 efficiency over transparency in these areas does not entail  
390 that we cease to value transparency here, and it certainly  
391 does not entail that decisions and actions in these areas are  
392 ethically neutral or trivial. It would be a mistake to regard  
393 areas in which our concern for efficiency wins out over our  
394 concern for transparency as areas that are “ethically neutral”,  
395 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 jeop-  
ard 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 goals<sup>8</sup> 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, stake-  
holders’ 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 medi-  
cal 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 pre-  
dictive 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 rela-  
tion 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 that a doc-  
tor’s extensive training results in sound medical judgment,  
and a well-functioning social system for verifying exper-  
tise (such as board certification and licensing). Trust is an  
essential feature of modern society’s successful (when it is

<sup>5</sup> See Skerker, Purves, and Jenkins (2015) on the anti-codifiability  
problem in robot and machine ethics.

<sup>6</sup> See London (2019) and Vincent (2018)

<sup>7</sup> See Robbins (2019) on valuing efficiency *rather than* transparency  
in certain non-trivial cases.

<sup>8</sup> 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.

442 successful) division of epistemic labor. It is clearly indis- 494  
 443 pensable for a successful doctor–patient relationship, and 495  
 444 the same holds for the epistemic and ethical relationships 496  
 445 between experts and laypeople in general. Trust is essential 497  
 446 in the absence of understanding and explanation (with suf- 498  
 447 ficient understanding and explanation, trust can be unnec- 499  
 448 essary). It is often thought that we trust processes that we 500  
 449 understand, as Riberio et al. (2016) make explicit here: 501

450 Whether humans are directly using machine learning 502  
 451 classifiers as tools, or are deploying models within 503  
 452 other products, a vital concern remains: if the users 504  
 453 do not trust a model or a prediction, they will not use 505  
 454 it. It is important to differentiate between two different 506  
 455 (but related) definitions of trust: (1) trusting a predic- 507  
 456 tion, i.e. whether a user trusts an individual prediction 508  
 457 sufficiently to take some action based on it, and (2) 509  
 458 trusting a model, i.e. whether the user trusts a model 510  
 459 to behave in reasonable ways if deployed. Both are 511  
 460 directly impacted by how much the human understands 512  
 461 a model's behaviour, as opposed to seeing it as a black 513  
 462 box. (Riberio, 2016, section 1, emphasis mine) 514

463 This is a common assumption regarding the relation 515  
 464 between trust and understanding, but it ignores the addi- 516  
 465 tional function and value of trust and trustworthiness men- 517  
 466 tioned above. Either increased trust or increased understand- 518  
 467 ing will typically result in an agent's increased willingness 519  
 468 to believe a certain decision is accurate or engage with a 520  
 469 certain tool. When patients trust their doctors, that trust is 521  
 470 not grounded in the patients' understanding of the doctors' 522  
 471 evidence or reasoning. This remains opaque. Patients trust 523  
 472 their doctors because they know that, in a well-functioning 524  
 473 social system which includes institutions dedicated to expert 525  
 474 verification, a person would not hold the position of doctor 526  
 475 unless they possessed the adequate expertise. 527

476 In a society that operates with a successful division of 528  
 477 epistemic labor, trust and trustworthiness can replace under- 529  
 478 standing as epistemically and ethically sound grounds for 530  
 479 belief. Laypeople believe the judgments of specialized 531  
 480 experts because they trust those experts—not because they 532  
 481 understand their reasoning—and they trust those experts 533  
 482 because their social framework includes institutions whose 534  
 483 role it is to verify the legitimacy of specialized experts. If 535  
 484 the ultimate aim of increased transparency is to establish 536  
 485 trustworthiness and build trust where appropriate, there may 537  
 486 be other avenues available for pursuing these goals—paths 538  
 487 that allow us to benefit from the power of opaque AI tech- 539  
 488 nology by verifying its trustworthiness. Transparency itself 540  
 489 need not be our goal. 541

490 If this is correct, then the options before us are either 542  
 491 (1) accept that the ethical concerns which give us reason to 543  
 492 employ opaque AI may outweigh the benefits of transpar- 544  
 493 ency, and determine how to best utilize opaque AI given

these epistemic limitations, or (2) refuse to employ opaque 494  
 AI in any ethically significant context on the grounds that 495  
 the use of an opaque process is ethically impermissible in 496  
 those contexts. 497

498 Given that there are enormous potential benefits that 499  
 could arise from the proper use of opaque AI in at least some 500  
 of the commonly identified ethically significant domains— 501  
 healthcare, climate science, the global food chain, public 502  
 safety—we would need powerful ethical reasons to sup- 503  
 port fully eliminating its use in these areas. The success of 504  
 the expert/layperson division of epistemic labor shows us 505  
 that many of our ordinary, ethically responsible, and reli- 506  
 able social practices already implicitly reject (2) above: we 507  
 routinely employ opaque processes in ethically significant 508  
 domains. And I will argue that there is no special reason to 509  
 embrace (2) in the case of AI while rejecting it in the case of 510  
 human experts. If this is correct, then we are left with option 511  
 (1), and the ethical question before us is no longer whether 512  
 we ought to allow opaque AI to operate in any ethically sig- 513  
 nificant domains but rather what are the most ethical ways 514  
 of harnessing opaque AI in these domains. 515

## 6 The expert/layperson relation—a blueprint for ethical opaque AI 516

517 I have suggested that we take our successful social practice 518  
 of deferring to specialized experts as a guide for develop- 519  
 ing an epistemically and ethically sound method for utiliz- 520  
 ing opaque AI. To this end, we will need to examine when 521  
 (i.e., under what conditions) it is epistemically and ethically 522  
 responsible to defer to experts rather than relying on one's 523  
 own reasoning. We also need to know what features make an 524  
 individual a genuine expert, how, as a society, we determine 525  
 that some individual is an expert, and what methods we use 526  
 for deciding how to act when multiple experts disagree in 527  
 their decisions. Fortunately, these questions are beginning to 528  
 receive increased attention both in sociology and philosophy, 529  
 under the general headings of social epistemology and the 530  
 epistemology of testimony.<sup>9</sup> 531

532 In what follows I will make a preliminary case for the 533  
 claim that the essential features of experts-- the features 534  
 that make expert opinion trustworthy, and our trust in those 535  
 individuals' decisions both epistemically and ethically 536  
 responsible—can be realized in AI as well. For this to be 537  
 the case, the relevant features of human experts must not be 538  
 essentially human features. Certainly, human experts have 539  
 noteworthy features that AI models lack; for instance, we 540  
 typically assume that human experts have a concept of the 541  
 greater good and a desire to promote it. If such traits play an 542  
 indispensable role in generating the trust and trustworthiness 543

<sup>9</sup> see Goldman 2001; Goldman 2014; Lackey 2016. 544

542 on which the expert/layperson relation depends, this relation  
 543 will not be a viable model for the ethical use of opaque AI.  
 544 As I hope to show below, the trust that exists in the expert/  
 545 layperson relation is not fundamentally based on faith in  
 546 the moral goodness of the expert but rather on the nature  
 547 of expertise and the existence of institutions that serve to  
 548 verify these experts. If these features are not uniquely human  
 549 features, then, insofar as we have ethically acceptable meth-  
 550 ods of evaluating when we ought to defer to human experts  
 551 in high-stakes contexts, we have a potential framework for  
 552 determining when it is ethically appropriate to defer to the  
 553 decisions generated by opaque AI technology.

554 In the mid 1980's, philosopher John Hardwig sparked  
 555 renewed interest in the social aspects of knowledge-building  
 556 by drawing attention to the myriad situations in which we  
 557 are better off—rationally speaking—deferring to someone  
 558 else's judgment on a particular matter rather than attempting  
 559 to reason through that matter ourselves. These are situations  
 560 in which the matter at hand concerns an area of highly spe-  
 561 cialized knowledge, and there are highly trained experts who  
 562 specialize in that area. In such a case, a layperson would be  
 563 more rationally justified in deferring to the expert's judg-  
 564 ment than they would in performing their own independ-  
 565 ent reasoning and standing by the judgment at which they  
 566 themselves had arrived. That is to say, a layperson has better  
 567 reasons to believe an expert's judgment is correct than his  
 568 or her own, even when that judgment conflicts with theirs.  
 569 Assuming that the layperson is a genuine layperson, and the  
 570 expert a genuine expert, Hardwig writes,

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

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

590 1. The laypeople have not, themselves, performed the rea-  
 591 soning that is being left to the expert.

2. The laypeople are not capable of performing the rea-  
 soning 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 Hard-  
 wig'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 tech-  
 nology 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 particu-  
 lar 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 applica-  
 tions, evaluating job candidates, deciding how to deploy  
 police resources, predicting effects of climate and weather  
 events on food production—have all previously been per-  
 formed 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 unprec-  
 edented. 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 ethi-  
 cally 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 follow-  
 ing section.)

2. Transparent models are not practically capable of per-  
 forming 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

639 researchers at the time the decision is being made. Rudin  
640 (2019) points to this aspect of the problem when she writes,

641 Black box models seem to uncover ‘hidden patterns’.  
642 The fact that many scientists have difficulty construct-  
643 ing interpretable models may be fueling the belief that  
644 black boxes have the ability to uncover subtle hidden  
645 patterns in the data about which the user was not pre-  
646 viously aware. A transparent model may be able to  
647 uncover these same patterns. If the pattern in the data  
648 was important enough that a black box model could  
649 leverage it to obtain better predictions, an interpretable  
650 model might also locate the same pattern and use it.  
651 Again, this depends on the ML researcher’s abil-  
652 ity to create accurate yet interpretable models. The  
653 researcher needs to create a model that has the capa-  
654 bility of uncovering the types of pattern that the user  
655 would find interpretable, but also the model needs to  
656 be flexible enough to fit the data accurately. This, and  
657 the optimization challenges discussed above, are where  
658 the difficulty lies with constructing interpretable mod-  
659 els. (2019, p. 201)

660 If equally proficient transparent models<sup>10</sup> already exist  
661 or could realistically be developed within the requisite  
662 timeframe (where ‘equally proficient’ takes into account  
663 the *speed* required to perform the task effectively as well  
664 as the scale of the task), the additional value conferred by  
665 their transparency may make them ethically preferable to  
666 an opaque model. Though Rudin is optimistic regarding the  
667 potential of transparent (in this case, interpretable) models to  
668 perform as well as opaque models, this is by no means guar-  
669 anteed. As she acknowledges, “This problem is compounded  
670 by the fact that researchers are now trained in deep learning,  
671 but not in interpretable ML...” and “It could be possible  
672 that there are application domains where a complete black  
673 box is required for a high stakes decision,” though she notes  
674 that, “As of yet, I have not encountered such an application”  
675 (2019, p. 207).

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

10FL01 <sup>10</sup> While “opaque” has a standard meaning in the literature on this  
10FL02 topic, “transparent” has several common meanings when used in  
10FL03 the context of AI models. A satisfactorily transparent AI model  
10FL04 might be an interpretable model, or an explicable model, or it may  
10FL05 be comprehensible to the relevant practitioner or stakeholder, etc. A  
10FL06 thorough account of how “transparency” has been interpreted in the  
10FL07 literature on AI regulations is beyond the scope of this discussion, but  
10FL08 see Chen 2018; Li et al. 2018; Lipton 2016; Miller 2017; Mittelstadt  
10FL09 et al. 2019; Molnar 2019; Riberio 2016; Rudin 2019; Zerilli 2002;

679 An explanation of an AI model would allow us to “assess  
680 the merits” of the evidence on which the model is basing its  
681 decision and “understand... how [the evidence] supports”  
682 that decision. The third criterion roughly specifies that in  
683 order for us to sacrifice transparency for the benefits gained  
684 by employing opaque AI in a particular ethically significant  
685 context, that opacity must be a result of our genuine inabil-  
686 ity to explain the operations of the AI model, rather than  
687 an unwillingness to deploy sufficient resources to the task.  
688 (Note that this issue will only arise when there is a question  
689 of irresponsibly employing opaque AI—the context itself  
690 must be ethically significant for ethical concerns to compete  
691 with the value of the efficiency or accuracy gained by utiliz-  
692 ing opaque AI models.)

693 In addition to this cursory description of when it would  
694 be reasonable for a layperson to defer to the judgment of an  
695 expert, Hardwig also provides a rough approximation of the  
696 personal features that make an individual an expert. Briefly,  
697 an expert must have engaged in “inquiry that has been sus-  
698 tained, prolonged, and systematic” (1985, p. 338). Though  
699 we would need to determine what features of an AI model  
700 would make its “inquiry” into a specific domain suitably  
701 “sustained, prolonged, and systematic,” this criterion seems  
702 to pose no special difficulty for AI. And given that these  
703 models fundamentally function by discovering and attuning  
704 themselves to patterns in data, such data-processing opera-  
705 tions 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)

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

724 In the familiar cases of human experts, the answers to  
725 both of these questions rely, in part, on the existence of  
726 a larger network of experts in addition to the individual  
727 (potential) expert in question. In areas of technical special-  
728 ization, (academic research, professions such as journalism  
729 and law, etc.) we commonly find established institutions and



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

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

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

771 Because of the importance that time plays in revealing  
 772 the reliability of expert decisions, of our method for veri-  
 773 fying individuals as genuine experts, and of our division  
 774 of epistemic labor in general, whatever way in which we  
 775 choose to adapt this feature to create an analogous method  
 776 for revealing the trustworthiness of any opaque AI technol-  
 777 ogy will be highly speculative. There are no obvious candi-  
 778 dates for artificial analogs of the passage of time. With that  
 779 in mind, one possible option would 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 the  
 782 history of AI models on a human timescale, it may be more

appropriate to frame the notion of “an adequate length of  
 time” on which to judge the reliability of an AI model to  
 reflect an AI timescale. So whereas ANNs and other deep  
 learning models may have emerged 10 years ago on a human  
 timescale, a massive number of epochs for those models has  
 passed within this span. (Determining the optimal number of  
 epochs for training a neural network is currently considered  
 something of an art in machine learning.)

783  
784  
785  
786  
787  
788  
789  
790  
791 There are a growing number of organizations dedicated  
 792 to developing something akin to a “credentialing processes”  
 793 for AI. The Institute of Electrical and Electronics Engineers  
 794 (IEEE) continuously updates its standards for the devel-  
 795 opment and use of AI. The International Organization for  
 796 Standardization (ISO) and The International Electrotechni-  
 797 cal Commission (IEC) both work to develop standards  
 798 that aim to make AI more “resilient, reliable, accurate, and  
 799 secure”. And the European Commission’s 2021 proposal for  
 800 Regulation on Artificial Intelligence includes a legal frame-  
 801 work by which to judge the risk of AI. The UK Institute for  
 802 Ethical AI and Machine Learning, the Global Partnership on  
 803 AI (GPAI), and the OECD AI Policy Observatory all support  
 804 projects and policy aimed at increasing trustworthiness in  
 805 AI. What form a successful credentialing process will eventu-  
 806 ally take, and to what extent these certification systems are  
 807 already in place, is a question to be addressed elsewhere. But  
 808 if we are interested in developing an approval process that  
 809 could certify certain opaque AI technology and approve its  
 810 use in particular contexts while allowing the technology to  
 811 remain opaque, we might make progress on this issue by  
 812 continuing research into the relevant features of familiar and  
 813 successful practices of certifying human experts.

814 The final feature of the expert-layperson relationship that  
 815 we will address here—our methods for dealing with cases of  
 816 expert disagreement—is simple to adapt in theory (though  
 817 perhaps less so in practice). Our successful division of epis-  
 818 temic labor crucially depends on the existence of multiple  
 819 independently trained experts in a single field, addressing  
 820 the same issue and converging on the same opinion through  
 821 a variety of independent methods. At the present moment,  
 822 it is unclear whether there exists a sufficient number—and  
 823 variety—of AI models that could perform the same ethically  
 824 significant task (whatever this task may be) to deal with dis-  
 825 agreement in an analogous way. But there may be no better  
 826 way to establish the requisite level of trustworthiness<sup>11</sup> [1]  
 827 of an opaque AI model than developing multiple, independ-  
 828 ent, opaque models, operating with distinct architecture and  
 829 trained on distinct (but appropriately relevant) data sets, and  
 830 finding that they converge on the same decision. Given that  
 831 opaque AI will be an ever-present ethical issue, developing

<sup>11</sup> 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.

832 multiple models to perform the same ethically significant  
833 task may be well worth the investment.

## 834 **7 Preliminary guidelines for the ethical use** 835 **of opaque AI**

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

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

851 (A) Is the context in question an ethically significant con-  
852 text? The process of evaluation begins with question (A): Is  
853 the context in question an ethically significant context? If we  
854 can be reasonably certain that the answer to (A) is “no,” then  
855 the ethical concerns surrounding the use of opaque AI do  
856 not arise in this situation, and we are at liberty to use opaque  
857 AI for the task at issue. Note that the triviality or ethical  
858 significance of a context will most often be decided accord-  
859 ing to a broad and diverse set of standards, some of which  
860 may involve apparently objective and quantifiable measures  
861 (for example, the potential consequences of utilizing some  
862 proposed AI technology in the global food supply chain) and  
863 some of which may involve standards that will vary relative  
864 to a cultural or social context (the impact of utilizing some  
865 proposed AI on the representation of a particular socially  
866 marginalized or vulnerable group). Note also that the ethical  
867 significance of a context will be a matter of degree, depend-  
868 ing on the gravity of the particular situation(s) involved. If  
869 the answer to (A) is “yes,” then we need to address question  
870 (B). Could the task at issue be performed equally well by a  
871 transparent process (whether human or AI)? This question  
872 will be familiar from the criteria for rationally deferring to  
873 experts in general. The additional benefits that arise from  
874 transparency in how decisions are made in all ethically sig-  
875 nificant contexts may outweigh whatever benefits the opaque  
876 AI may provide. Here, it is important to note that “per-  
877 forming a task equally as well” will include—at minimum—  
878 issues of equity and fairness in addition to efficiency and  
879 accuracy. As noted in Sect. 4 3, we cannot entirely ignore

880 the harms of opportunity costs for the sake of eliminating  
881 bias, especially when those costs are borne by marginalized  
882 and vulnerable populations. To permit the use of opaque AI  
883 in an ethically significant context, the answer to question  
884 (A) must be yes, and the answer to question (B) must be no.  
885 If so, then it may be ethically permissible to utilize opaque  
886 AI, subject to further consideration, such as those raised in  
887 question (C). Are the benefits of successfully performing  
888 this task greater than both (i) the cost of potentially failing  
889 at this task (whatever constitutes “failure” in this case) and  
890 (ii) the cost of not performing this task at all? If there are  
891 ethically significant cases in which all three bars are met,  
892 then there are non-trivial cases in which we would be per-  
893 mitted—perhaps required—to utilize opaque AI. And given  
894 that meeting all three bars requires that the opaque model in  
895 question be reliable and trustworthy, we will need a frame-  
896 work for evaluating the reliability and trustworthiness of  
897 opaque AI technology. I hope to have made a preliminary  
898 case for looking to our successful social practice of deferring  
899 to experts in ethically significant domains for a blueprint of  
900 how to responsibly employ opaque AI in such a case.

## 901 **8 Conclusion**

902 I acknowledge that, even as guidelines go, those given above  
903 are considerably vague. I view this vagueness as appropri-  
904 ate, and—practically speaking—ineliminable. Here, we are  
905 concerned with developing rules for ethical action in the use  
906 of AI, and as Aristotle said, we should only look for preci-  
907 sion in each class of things just so far as the nature of the  
908 subject admits. Any rule, no matter how precise, requires  
909 interpretation when applied to a particular case. And when  
910 the interpretation of those rules involves disentangling and  
911 weighing competing moral values, it is the process of inter-  
912 pretation itself—and not the rule—that will be doing the  
913 lion’s share of the work. So I would suggest that insofar  
914 as these guidelines are vague, their vagueness is appropri-  
915 ate to the subject at hand. Deciding whether a task could  
916 be performed equally well by some satisfactorily transpar-  
917 ent (human or algorithmic) decision-making process will  
918 involve weighing competing values, and the relative strength  
919 of those competing values will depend on the ethical inclina-  
920 tions of the individuals performing the evaluation. There is  
921 no standard, universally applicable measure for assigning  
922 weights to these values; each case will need to be evalu-  
923 ated individually, and an argument will need to be made for  
924 weighting any of these values more strongly than the others.  
925 The same is true for deciding whether the benefits of success  
926 are worth the potential costs of failure. Human judgment  
927 cannot be entirely removed from decision-making in ethi-  
928 cally significant domains, no matter how trustworthy the AI  
929 technology involved. At minimum, humans must still be

930 in-the-loop to (1) make case-specific value-judgments, and  
 931 (2) make cost/benefit assessments in cases where the costs  
 932 and benefits are not fully commensurable. And given that  
 933 we are discussing opaque AI technology, humans will need  
 934 to be in-the-loop to monitor for potential instances of biased  
 935 outcomes. The threat of bias will remain, whether or not  
 936 the cost of that potential bias is outweighed by the potential  
 937 benefits of a successful outcome.

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

955 **Funding** No funding was received to assist with the preparation of  
 956 this manuscript.

## 957 Declarations

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

## 960 References

- 961 Ahmed M (2018) Aided by Palantir, the LAPD uses predictive polic-  
 962 ing to monitor specific people and neighborhoods. *The Intercept*.  
 963 <https://theintercept.com/2018/05/11/predictive-policing-surveillance-los-angeles/>  
 964  
 965 Barocas S (2018) Accounting for artificial intelligence: rules, reasons,  
 966 rationales. In: Human rights, ethics, and artificial intelligence,  
 967 30 Nov. Harvard Kennedy School Carr Center for Human Rights  
 968 Policy. Lecture  
 969 Barry-Jester A, Casselman B, Goldstein D (2015) The new science of  
 970 sentencing. *The Marshall Project*. <https://www.themarshallproject.org/2015/08/04/the-new-science-of-sentencing>  
 971  
 972 Berk RA, Sorenson SB, Barnes G (2016) Forecasting domestic vio-  
 973 lence: a machine learning approach to help inform arraignment  
 974 decisions. *J Empir Leg Stud* 13(1):94–115. <https://doi.org/10.1111/jels.12098>  
 975  
 976 Chen C et al (2018) This looks like that: deep learning for interpretable  
 977 image recognition. Preprint at <https://arxiv.org/abs/1806.10574>  
 978 de Bruijne M (2016) Machine learning approaches in medical image  
 979 analysis: from detection to diagnosis. *Med Image Anal* 33:94–97.  
 980 <https://doi.org/10.1016/j.media.2016.06.032>

- Dhar J, Ranganathan A (2015) Machine learning capabilities in medi-  
 cal diagnosis applications: computational results for hepatitis dis-  
 ease. *Int J Biomed Eng Technol* 17(4):330–340. <https://doi.org/10.1504/IJBET.2015.069398>  
 Erickson BJ, Korfiatis P, Akkus Z, Kline TL (2017) Machine learning  
 for medical imaging. *Radiographics* 37(2):505–515. <https://doi.org/10.1148/rg.2017160130>  
 Ensign, D., Friedler, S. A., Neville, S., Scheidegger, C., & Venkatasu-  
 bramanian, S. (2017). Run away feedback loops in predictive  
 policing. In *Proceedings of machine learning research*, 81, 1–12.  
 Retrieved from <http://arxiv.org/abs/1706.09847>  
 European Commission (2019) Ethics Guidelines for Trustworthy AI,  
<https://ec.europa.eu/futurium/en/ai-allianceconsultation>.  
 Floridi L, Cowls J, Beltrametti M, Chatila R, Chazerand P, Dignum  
 V, Vayena E (2018) AI4People—an ethical framework for a  
 good AI society: opportunities, risks, principles, and recom-  
 mendations. *Mind Mach* 28(4):689–707. <https://doi.org/10.1007/s11023-018-9482-5>  
 Future of Life Institute (2017) Asilomar AI Principles. <https://futureoflife.org/ai-principles/>  
 Goldman AI (2001) Experts: which ones should you trust? *Philos Phenomenol Res* 63(1):85–110  
 Goldman AI (2014) Social process reliabilism: solving justification  
 problems in collective epistemology. *Lackey* 2014:11–41. <https://doi.org/10.1093/acprof:oso/9780199665792.003.0002>  
 Günther M, Kasirzadeh A (2022) Algorithmic and human decision  
 making: for a double standard of transparency. *AI & Soc*. <https://doi.org/10.1007/s00146-021-01200-5>  
 Hardwig J (1985) Epistemic dependence. *J Philos* 82(7):335–349  
 Joh, E. E. (2017). Feeding the machine: Policing, crime data, & algo-  
 rithms. *William & Mary Bill of Rights Journal*, 26, 287.  
 Lackey J (2016) What is justified group belief? *Philos Rev* 125(3):341–  
 396. <https://doi.org/10.1215/00318108-3516946>  
 Li O, Liu H, Chen C, Rudin C (2018) Deep learning for case-based  
 reasoning through prototypes: a neural network that explains its  
 predictions. In: *Proceedings of AAAI Conference on Artificial  
 Intelligence* 3530–3537 (AAAI, 2018).  
 Lipton ZC (2016) The mythos of model interpretability. In: *ICML  
 Workshop on Human Interpretability in Machine Learning*, vol  
 2017, pp. 96–100, 24  
 London AJ (2019) Artificial intelligence and black-box medical deci-  
 sions: accuracy versus explainability. *Hastings Center Rep* 49.  
<https://doi.org/10.1002/hast.973>  
 Miller T (2017) Explanation in artificial intelligence: insights from the  
 social sciences. arXiv  
 Mittelstadt B, Russell C, Wachter S (2019) Explaining explanations in  
 AI. In: *Proceedings of fairness, accountability, and transparency  
 (FAT\*)* (ACM, 2019)  
 Molnar C (2019) Interpretable machine learning  
 Nadella S (2016) Microsoft’s CEO explores how humans and AI Can  
 solve society’s challenges—together. *Slate*. <https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societyschallenges.html>  
 O’Neil C (2016) Weapons of math destruction: how big data increases  
 inequality and threatens democracy. Crown  
 Ribeiro MT, Singh S, Guestrin C (2016) “Why should I trust you?”:  
 explaining the predictions of any classifier. *KDD*  
 Robbins S (2019) A misdirected principle with a catch: explicabil-  
 ity for AI. *Mind Mach* 29:495–514. <https://doi.org/10.1007/s11023-019-09509-3>  
 Rudin C (2019) Stop explaining black box machine learning models  
 for high stakes decisions and use interpretable models instead. *Nat  
 Mach Intell* 1:206–215  
 Skerker M, Purves D, Jenkins R (2015) Autonomous machines, moral  
 judgment, and acting for the right reasons. *Ethical Theory Moral  
 Pract* 18(4):851–872 (**Special Issue: BSET-2014**)

- 1047 Vincent J (2018) AI that detects cardiac arrests during emergency calls will be tested across Europe this summer. The  
1048 Verge. [https://www.theverge.com/2018/4/25/17278994/](https://www.theverge.com/2018/4/25/17278994/ai-cardiac-arrest-corti-emergency-call-response)  
1049 ai-cardiac-arrest-corti-emergency-call-response  
1050  
1051 Zerilli J (2022) Explaining machine learning decisions. *Philos Sci*  
1052 89(1):1–19  
1053 Zerilli J, Knott A, Maclaurin J et al (2019) Transparency in algorithmic  
1054 and human decision-making: is there a double standard? *Philos*  
1055 *Technol* 32:661–683. <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