

1 **Draft NIST Special Publication 1270**

2
3
4 **A Proposal for Identifying and**
5 **Managing Bias in Artificial**
6 **Intelligence**

7
8
9 Reva Schwartz
10 Leann Down
11 Adam Jonas
12 Elham Tabassi

13
14 This draft publication is available free of charge from:
15 <https://doi.org/10.6028/NIST.SP.1270-draft>

Draft NIST Special Publication 1270

A Proposal for Identifying and Managing Bias within Artificial Intelligence

Reva Schwartz

*National Institute of Standards and Technology
Information Technology Laboratory*

Leann Down

Adam Jonas
Parenthetic, LLC

Elham Tabassi

*National Institute of Standards and Technology
Information Technology Laboratory*

This draft publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.1270-draft>

June 2021



U.S. Department of Commerce
Gina M. Raimondo, Secretary

National Institute of Standards and Technology
*James K. Olthoff, Performing the Non-Exclusive Functions and Duties of the Under Secretary of Commerce
for Standards and Technology & Director, National Institute of Standards and Technology*

65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

Certain commercial entities, equipment, or materials may be identified in this document in order to describe an experimental procedure or concept adequately. Such identification is not intended to imply recommendation or endorsement by the National Institute of Standards and Technology, nor is it intended to imply that the entities, materials, or equipment are necessarily the best available for the purpose.

**National Institute of Standards and Technology Special Publication 1270 (Draft)
Natl. Inst. Stand. Technol. Spec. Publ. 1270 (Draft),
30 pages (June 2021)**

**This draft publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.1270-draft>**

Organizations are encouraged to review this draft publication during the public comment period and provide feedback to NIST.

Public comment period: June 21 – September 10, 2021
National Institute of Standards and Technology
Attn: Information Technology Laboratory
100 Bureau Drive
Gaithersburg, Maryland 20899-2000 70 Email: ai-bias@list.nist.gov

101 **Abstract**

102 NIST contributes to the research, standards, evaluation, and data required to advance the
103 development and use of trustworthy artificial intelligence (AI) to address economic, social, and
104 national security challenges and opportunities. Working with the AI community, NIST has
105 identified the following technical characteristics needed to cultivate trust in AI systems:
106 accuracy, explainability and interpretability, privacy, reliability, robustness, safety, and security
107 (resilience) – and that harmful biases are mitigated. Mitigation of risk derived from bias in AI-
108 based products and systems is a critical but still insufficiently defined building block of
109 trustworthiness. This report proposes a strategy for managing AI bias, and describes types of bias
110 that may be found in AI technologies and systems. The proposal is intended as a step towards
111 consensus **standards** and a **risk-based framework** for trustworthy and responsible AI. The
112 document, which also contains an alphabetical glossary that defines commonly occurring biases
113 in AI, contributes to a fuller description and understanding of the challenge of harmful bias and
114 ways to manage its presence in AI systems.

115

116

117 **Key words**

118 bias, trustworthiness, AI safety, AI lifecycle, AI development

119

120	Table of Contents	
121	TABLE OF CONTENTS	ii
122	1. INTRODUCTION	1
123	2. THE CHALLENGE POSED BY BIAS IN AI SYSTEMS	2
124	3. APPROACH	4
125	4. IDENTIFYING AND MANAGING BIAS IN ARTIFICIAL INTELLIGENCE	5
126	<i>Figure 1: A three-stage approach for managing AI bias</i>	6
127	PRE-DESIGN STAGE	7
128	PROBLEM FORMULATION AND DECISION MAKING	7
129	OPERATIONAL SETTINGS AND UNKNOWN IMPACTS	7
130	OVERSELLING TOOL CAPABILITIES AND PERFORMANCE	7
131	PRACTICES	8
132	REAL-WORLD EXAMPLE	8
133	DESIGN AND DEVELOPMENT STAGE	8
134	OPTIMIZATION OVER CONTEXT	9
135	PRACTICES	9
136	REAL-WORLD EXAMPLE	10
137	DEPLOYMENT STAGE	10
138	DISCRIMINATORY IMPACT	10
139	INTENDED CONTEXT VS ACTUAL CONTEXT	10
140	CONTEXTUAL GAPS LEAD TO PERFORMANCE GAPS	11
141	PRACTICAL IMPROVEMENTS	12
142	<i>Figure 2: Example of bias presentation in three stages modeled on the AI lifecycle.</i>	12
143	5. CONCLUSION AND NEXT STEPS	13
144	6. APPENDICES	14
145	APPENDIX A: GLOSSARY	14
146	APPENDIX B: COLLABORATIVE WORK	17
147	7. REFERENCES	18
148		
149		

150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184

Acknowledgments

The authors wish to thank the many people who assisted with the development of this document, including our NIST colleagues, and the many academic and technical reviewers who took the time to provide their valuable feedback.

Audience

The main audience for this document is researchers and practitioners in the field of trustworthy and responsible artificial intelligence. Researchers will find this document useful for understanding a view of the challenge of bias in AI, and as an initial step toward the development of standards and a risk framework for building and using trustworthy AI systems. Practitioners will benefit by gaining an understanding about bias in the use of AI systems.

Trademark Information

All trademarks and registered trademarks belong to their respective organizations.

Note to Reviewers

As described throughout this report, one goal for NIST’s work in trustworthy AI is the development of a **risk management framework** and accompanying **standards**. To make the necessary progress towards that goal, NIST intends to carry out a variety of activities in 2021 and 2022 in each area of the core building blocks of trustworthy AI (accuracy, explainability and interpretability, privacy, reliability, robustness, safety, and security (resilience), and mitigation of harmful bias). This will require a concerted effort, drawing upon experts from within NIST and external stakeholders. NIST seeks additional collaborative feedback from members of the research, industry, and practitioner community throughout this process. All interested parties are encouraged to please submit comments about this draft report, and the types of activities and events which would be helpful, via the public comment process described on page 3 of this document. There will also be opportunities for engaging in discussions about and contributing to development of key practices and tools to manage Bias in AI. Please look for announcements for webinars, call for position papers, and request for comment on NIST document(s).

185 1. Introduction

186 The National Institute of Standards and Technology (NIST) promotes U.S. innovation and
187 industrial competitiveness by advancing measurement science, standards, and technology in
188 ways that enhance economic security and improve our quality of life. Among its broad range of
189 activities, NIST contributes to the research, standards, evaluations, and data required to advance
190 the development, use, and assurance of trustworthy artificial intelligence (AI).

191
192 In August 2019, fulfilling an assignment in an Executive Order¹ on AI, NIST released “A Plan
193 for Federal Engagement in Developing Technical Standards and Related Tools.” [100] Based on
194 broad public and private sector input, this plan recommended a deeper, more consistent, and
195 long-term engagement in AI standards “to help the United States to speed the pace of reliable,
196 robust, and trustworthy AI technology development.” NIST research in AI continues along this
197 path to focus on how to measure and enhance the trustworthiness of AI systems. Working with
198 the AI community, NIST has identified the following technical characteristics needed to cultivate
199 trust in AI systems: accuracy, explainability and interpretability, privacy, reliability, robustness,
200 safety, and security (resilience) – and that harmful biases are mitigated.

201
202 This paper, *A Proposal for Identifying and Managing Bias in Artificial Intelligence*, has been
203 developed to advance methods to understand and reduce harmful forms of AI bias. It is one of a
204 series of documents and workshops in the pursuit of a **framework for trustworthy and**
205 **responsible AI**.

206
207 While AI has significant potential as a transformative technology, it also poses inherent risks.
208 One of those risks is bias. Specifically, how the presence of bias in automated systems can
209 contribute to harmful outcomes and a public lack of trust. Managing bias is a critical but still
210 insufficiently developed building block of trustworthiness.

211
212 The International Organization for Standardization (ISO) defines bias in statistical terms: “the
213 degree to which a reference value deviates from the truth” [67]. This deviation from the truth can
214 be either positive or negative, it can contribute to harmful or discriminatory outcomes or it can
215 even be beneficial. From a societal perspective, bias is often connected to values and viewed
216 through the dual lens of differential treatment or disparate impact, key legal terms related to
217 direct and indirect discrimination, respectively.

218
219 Not all types of bias are negative, and there many ways to categorize or manage bias; this report
220 focuses on biases present in AI systems that can lead to harmful societal outcomes. These
221 harmful biases affect people’s lives in a variety of settings by causing disparate impact, and
222 discriminatory or unjust outcomes. The presumption is that bias is present throughout AI
223 systems, the challenge is identifying, measuring, and managing it. Current approaches tend to
224 classify bias by type (i.e.: statistical, cognitive), or use case and industrial sector (i.e.: hiring,
225 health care, etc.), and may not be able to provide the broad perspective required for effectively
226 managing bias as the context-specific phenomenon it is. This document attempts to bridge that

¹<https://www.federalregister.gov/documents/2019/02/14/2019-02544/maintaining-american-leadership-in-artificial-intelligence>

227 gap and proposes an approach for managing and reducing the impacts of harmful biases² *across*
228 contexts. The intention is to leverage key locations within stages of the AI lifecycle for optimally
229 identifying and managing bias. As NIST develops a framework and standards in this area, the
230 proposed approach is a starting point for community-based feedback and follow-on activities
231 related to bias and its role in trustworthy AI.
232

233 2. The Challenge Posed by Bias in AI Systems

234 The proliferation of modeling and predictive approaches based on data-driven and machine
235 learning techniques has helped to expose various social biases baked into real-world systems,
236 and there is increasing evidence that the general public has concerns about the risks of AI to
237 society. Distrust in AI can manifest itself through a belief that biases may be automated within
238 these technologies, and can perpetuate harms more quickly, extensively, and systematically than
239 human and societal biases on their own. Human decisions based on automated and predictive
240 technology are often made in settings such as hiring or criminal justice, and can create harmful
241 impacts and amplify and accelerate existing social inequities or, at minimum, perceptions of
242 inequities. While it's unlikely that technology exhibiting "zero risk" can be developed, managing
243 and reducing the impacts of harmful biases in AI is possible and necessary.
244

245 Public attitudes about AI technology suggest that, while often depending on the application, most
246 Americans are unaware when they are interacting with AI enabled tech [53] but feel there needs
247 to be a "higher ethical standard" than with other forms of technologies [76]. This mainly stems
248 from the perceptions of fear of loss of control and privacy [47,125,133,137]. Certainly, there is
249 no shortage of examples where bias in some aspect of AI technology and its use has caused harm
250 and negatively impacted people's lives, such as in hiring [5,12,16,17,36,62,118], health care
251 [46,52,55,59,83,88,103,122,123], and criminal justice [7,20,29,41,44,56,66,74,75,78,87,
252 140,142]. Indeed, there are many instances in which the deployment of AI technologies have
253 been accompanied by concerns of whether and how societal biases are being perpetuated or
254 amplified [3,10,14,15,22,24,34,42,45,61,102,105,108,116,126,139].
255

256 Since AI systems are deployed across various contexts, the associated biases that come with their
257 use create harm in context-specific ways. This proliferation of AI bias into an ever-increasing list
258 of settings makes it especially difficult to develop overarching guidance or mitigation
259 techniques. A confounding factor is that it is especially difficult to predict where and how AI
260 systems will be used. A current approach to the challenge of AI bias is to tackle a given use case
261 where a particularly prevalent type of bias resides. This ad-hoc strategy is difficult to scale, and
262 is unlikely to achieve what is required for building systems that the public can trust. Instead of
263 viewing the challenge of AI bias within a given context or use case, a broader perspective can
264 strike the problem of AI bias where it might be easiest to manage – within the design,
265 development, and use of AI systems.
266

267 There are specific conditional traits associated with automation that exacerbate distrust in AI
268 tools. One major purpose, and a significant benefit, of automated technology is that it can make
269 sense of information more quickly and consistently than humans. There have long been two
270 common assumptions about the rise and use of automation: it could make life easier [137] and

² For the purpose of this document the term "managing bias" will be used to refer to approaches for managing, reducing or mitigating bias.

271 also create conditions that reduce (or eliminate) biased *human* decision making and bring about a
272 more equitable society [78]. These two tenets have led to the deployment of automated and
273 predictive tools within trusted institutions and high-stake settings. While AI can help society
274 achieve significant benefits, the convenience of automated classification and discovery within
275 large datasets may come with a potentially significant downside. As these tools proliferate across
276 our social systems, there has been increased interest in identifying and mitigating their harmful
277 impacts.

278
279 The difficulty in characterizing and managing AI bias is exemplified by systems built to model
280 concepts that are only partially observable or capturable by data. Without direct measures for
281 these often highly complex considerations, AI development teams often use proxies. For
282 example, for “criminality,” a measurable index, or construct, might be created from other
283 information, such as past arrests, age, and region. For “employment suitability,” an AI algorithm
284 might rely on time in prior employment, previous pay levels, education level, participation in
285 certain sports [115], or distance from the employment site [51] (which might disadvantage
286 candidates from certain neighborhoods).

287
288 There are many challenges that come with this common practice (see [89] for a thorough
289 review). One challenge rests on the reality that decisions about which data to use for these
290 indices are often made based on what is available or accessible, rather than what might be most
291 suitable - but difficult or impossible to utilize [49]. Relatedly, instead of identifying specific
292 questions of interest *first*, researchers, developers, and practitioners may “go where the data is”
293 and adapt their questions accordingly [130]. Data can also differ significantly between what is
294 collected and what occurs in the real world [71,72,109]. For example, responses to online
295 questionnaires are from a specific sampling of the kinds of people who are online, and therefore
296 leaves out many other groups. Data representing certain societal groups may be excluded in the
297 training datasets used by machine learning applications [40]. And, datasets used in natural
298 language processing often differ significantly from their real-world applications [113] which can
299 lead to discrimination [128] and systematic gaps in performance.

300
301 Even if datasets are reflective of the real world, they may still exhibit entrenched historical and
302 societal biases, or improperly utilize protected attributes. (Federal laws and regulations have
303 been established to prohibit discrimination based on grounds such as gender, age, and religion.)
304 Simply excluding these explicit types of attributes will not remedy the problem, however, since
305 they can be inadvertently inferred in other ways (for example, browsing history), and still
306 produce negative outcomes for individuals or classes of individuals [12]. So, the proxies used in
307 development may be both a poor fit for the concept or characteristic seeking to be measured, and
308 reveal unintended information about persons and groups.

309
310 Additionally, for much of the public, AI is not necessarily something with which they directly
311 interact, and systems' algorithmic assumptions may not be transparent to them. Nevertheless,
312 many people are affected or used as inputs by AI technologies and systems. This can happen
313 when an individual applies for a loan [136], college [48], or a new apartment [77]. Historical,
314 training data, and measurement biases are “baked-in” to the data used in the algorithmic models
315 underlying those types of decisions. Such biases may produce unjust outcomes for racial and

316 ethnic minorities in areas such as criminal justice [7,41,56,74,75,78,87,140,142], hiring
317 [4,5,12,16,17,36,118,119], and financial decisions [13,65].

318

319 Another cause for distrust may be due to an entire class of untested and/or unreliable algorithms
320 deployed in decision-based settings. Often a technology is not tested – or not tested extensively –
321 before deployment, and instead deployment may be used *as* testing for the technology. An
322 example is the rush to deploy systems during the COVID pandemic that have turned out to be
323 methodologically flawed and biased [117,124,141]. There are also examples from the literature
324 which describe technology that is based on questionable concepts, deceptive or unproven
325 practices, or lacking theoretical underpinnings [2,9,13,30,33,62,129,141]. The broad consensus
326 of the literature is that systems meant for decision making or predictive scenarios should
327 demonstrate validity and reliability under the very specific setting in which it is intended to be
328 deployed (hiring purposes, risk assessments in the criminal justice system, etc.). The decisions
329 based on these algorithms affect people’s lives in significant ways, and it is appropriate to expect
330 protections in place to safeguard from certain systems and practices. The public’s cautious
331 opinions toward AI [138] might turn increasingly negative if *new* technologies appear which are
332 based on the same approaches that have already contributed to systematic and well-documented
333 societal harms.

334

335 To summarize the problem, there are many reasons for potential public distrust of AI related to
336 bias in systems. These include:

337

- 338 • The use of datasets and/or practices that are inherently biased and historically contribute
339 to negative impacts
- 340 • Automation based on these biases placed in settings that can affect people’s lives, with
341 little to no testing or gatekeeping
- 342 • Deployment of technology that is either not fully tested, potentially oversold, or based on
343 questionable or non-existent science causing harmful and biased outcomes

344

345 Identifying and working to manage these kinds of bias can mitigate concerns about
346 trustworthiness for in-place and in-development AI technologies and systems. An effective
347 approach will likely need to be one that is not segmented by use case, but works across contexts.

348

349 Improving trust in AI systems can be advanced by putting mechanisms in place to reduce
350 harmful bias in both deployed systems *and* in-production technology. Such mechanisms will
351 require features such as a common vocabulary, clear and specific principles and governance
352 approaches, and strategies for assurance. For the most part, the standards for these mechanisms
353 and associated performance measurements still need to be created or adapted. The goal is not
354 “zero risk,” but to manage and reduce bias in a way that contributes to more equitable outcomes
355 that engender public trust. These challenges are intertwined in complex ways and are unlikely to
356 be addressed with a singular focus on one factor or within a specific use or industry.

357

358 3. Approach

In the lead-up to this report, the authors sought to capture common themes about the many ways
bias is defined and categorized in AI technology. This was accomplished through a literature

359 review, discussions with leaders in the field, a NIST-hosted workshop on bias in AI³, and the
360 evaluation of prominent topics across the broader AI research community. This work is not
361 without precedent; there are previous attempts to define and classify AI bias
362 [26,35,64,68,69,91,94,95,98,106,127].

363
364 The literature review consisted of a total of 313 articles, books, reports, and news publications
365 about AI bias⁴ from a variety of perspectives. In the survey of the literature, we identified a list
366 of prominent biases present in AI that are contributors to societal harms. This list and
367 accompanying definitions are presented in an alphabetical glossary in Appendix A.

368
369 The reviewed literature suggests that the expansion of AI into many aspects of public life
370 requires extending our view from a mainly technical perspective to one that considers AI within
371 the social system it operates [3,18,19,31,34,40,41,43,71,97,118,120,134]. Taking social factors
372 into consideration is necessary for achieving trustworthy AI, and can enable a broader
373 understanding of AI impacts and the key decisions that happen throughout, and beyond, the AI
374 lifecycle – such as whether technology is even a solution to a given task or problem [11,49].
375 Such a change in perspective will require working with new stakeholders and developing
376 guidance for effectively engaging social factors within a technical perspective. A key factor in
377 this area is the many ways in which institutions indirectly drive the design and use of AI. Also,
378 while AI practices may not intend to contribute to inequality or other negative forms of bias,
379 there are always complex social factors that may be overlooked, especially since biases play out
380 in context-specific ways and may not be captured or understood within one setting.

381
382 Whether statistical or societal, bias continues to be a challenge for researchers and technology
383 developers seeking to develop and deploy trustworthy AI applications. How bias and trust
384 interrelate is a key societal question, and understanding it will be paramount to improving
385 acceptance of AI systems. A consistent finding in the literature is the notion that trust can
386 improve if the public is able to interrogate systems and engage with them in a more transparent
387 manner. Yet, in their article on public trust in AI, Knowles and Richards state “...members of the
388 public do not need to trust individual AIs at all; what they need instead is the sanction of
389 authority provided by suitably expert auditors that AI can be trusted” [80]. Creating such an
390 authority requires standard practices, metrics, and norms. NIST has experience in creating
391 standards and databases, and has been evaluating the algorithms used in biometric technologies
392 since the 1960s. With the development of privacy and cybersecurity frameworks [99,101], NIST
393 has helped organizations manage risks of the digital environment, and, through a series of reports
394 and workshops, intends to contribute to a similar collaborative approach for managing AI
395 trustworthiness as part of broader stakeholder efforts.

396
397 4. Identifying and Managing Bias in Artificial Intelligence
398 Improving trust in AI by mitigating and managing bias starts with identifying a structure for how
399 it presents within AI systems and uses. We propose a three-stage approach derived from the AI

³ For more information about this workshop see <https://www.nist.gov/news-events/events/2020/08/bias-ai-workshop>, and Appendix B of this document.

⁴ The full bibliographic survey can be found at https://www.nist.gov/system/files/documents/2021/03/26/20210317_NIST%20AI_Bibliography.pdf

400 lifecycle, to enable AI designers and deployers to better relate specific lifecycle processes with
401 the types of AI bias, and facilitate more effective management of it. Organizations that design
402 and develop AI technology use the AI lifecycle to keep track of their processes and ensure
403 delivery of high-performing functional tools - but not necessarily to identify harms or manage
404 them. Currently, there is no single global or industrial AI lifecycle standard, but many versions
405 used across multiple sectors and regions with a range of stages. The approach for identifying and
406 managing AI bias proposed in this report is adapted from current versions of the AI lifecycle⁵,
407 and consists of three distinct stages, and presumed accompanying stakeholder groups. This
408 approach is a starting point and NIST seeks feedback about its viability and implementation.

- 409
- 410 1. **PRE-DESIGN**: where the technology is devised, defined and elaborated
 - 411 2. **DESIGN AND DEVELOPMENT**: where the technology is constructed
 - 412 3. **DEPLOYMENT**: where technology is used by, or applied to, various individuals or
413 groups.
- 414

415 Figure 1: A three-stage approach for managing AI bias



416
417
418 The following sub-sections provide key considerations and examples that highlight how
419 statistical biases present across various stages of AI applications; and reflect and interact with the
420 many human cognitive and societal biases that are inherent in the data, modeling, decision
421 making, and practical processes associated with the use of AI systems across sectors and
422 contexts.

423

⁵ The following AI lifecycles were utilized as key guidance for this report: Centers of Excellence (CoE) at the US General Services Administration [70] [IT Modernization CoE. (n.d.)], the Organisation for Economic Co-operation and Development [106] [Organisation for Economic Co-operation and Development. (2019)]. Another model of the AI lifecycle is currently under development with the Joint Technical Committee of the International Organization for Standardization (ISO) and the International Electrotechnical Commission (IEC) (see <https://www.iso.org/standard/81118.html>)

424 PRE-DESIGN STAGE

425

426 Problem formulation and decision making

427 AI products start in the pre-design stage, where planning, problem specification, background
428 research, and identification and quantification of data take place. Decisions here include how to
429 frame the problem, the purpose of the AI component, and the general notion that there is a
430 problem requiring or benefitting from a technology solution. Since many of the downstream
431 processes hinge on decisions from this stage, there is a lot of pressure here to “get things right.”
432 Central to these decisions is *who* (individuals or groups) makes them and which individuals or
433 teams have the most power or control over them. These early decisions and who makes them can
434 reflect individual and group heuristics and limited points of view, affect later stages and
435 decisions in complex ways, and lead to biased outcomes [12,31,43,72,109,120]. This is a key
436 juncture where well-developed guidance, assurance, and governance processes can assist
437 business units and data scientists to collaboratively integrate processes that reduce bias without
438 being cumbersome or blocking progress.

439

440 Operational settings and unknown impacts

441 Current assumptions in AI development often revolve around the idea of technological
442 solutionism – the perception that technology will lead to only *positive* solutions. This perception,
443 often combined with a singular focus on tool optimization, can be at odds with operational
444 scenarios, increasing the difficulty for the practitioners who have to make sense of tool output –
445 often in high stakes settings [96]. What seems like a good idea for how a given dataset can be
446 utilized in a specific use case might be perceived differently by the systems’ end users or those
447 affected by the systems’ decisions. It is an obvious risk to build algorithmic-based decision tools
448 for settings already known to be discriminatory. Yet, awareness of which conditions will lead to
449 disparate impact or other negative outcomes is not always apparent in pre-design, and can be
450 easily overlooked once in production.

451

452 Overselling tool capabilities and performance

453 Whether unconscious or unintentional, pre-design is often where decisions are made that can
454 inadvertently lead to harmful impact, or be employed to extremely negative societal ends. By not
455 addressing the possibility of optimistic and potentially inflated expectations related to AI
456 systems, risk management processes could fail to communicate and set reasonable limits related
457 to mitigating such potential harms. In extreme cases, with tools or apps that are fraudulent,
458 pseudoscientific, prey on the user, or generally exaggerate claims, the goal should not be to
459 ensure tools are bias-free, but to reject the development outright in order to prevent
460 disappointment or harm to the user as well as to the reputation of the provider.

461

462 Other problems that can occur in pre-design include poor problem framing, basing technology on
463 spurious correlations from data-driven approaches, failing to establish appropriate underlying
464 causal mechanisms, or generally technically flawed [22,34,40,52,54,89,102,110]. In such cases
465 (often termed “fire, ready, aim”), the solution may not be mitigation, but rather, rejection of the
466 system or the way in which the perceived underlying problem is framed. These types of
467 scenarios may reinforce public distrust of AI technology as systems that are untested or
468 technically flawed can also contribute to bias. Technology designed for use in high-stakes
469 settings requires extensive testing to demonstrate valid and reliable performance [58,112].

470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515

Practices

There is currently momentum for AI researchers to include statements about the potential societal impacts [114] when submitting their work to journals or conferences. Identifying and addressing potential biases early in the problem formulation process is an important step in this process. It is also complicated by the role of power and decision making [96]. A consistent theme from the literature is the benefit of engaging a variety of stakeholders and maintaining diversity along social lines where bias is a concern (racial diversity, gender diversity, age diversity, diversity of physical ability) [32]. These kinds of practices can lead to a more thorough evaluation of the broad societal impacts of technology-based tools across the three stages. Identifying downstream impacts may take time and require the involvement of end-users, practitioners, subject matter experts, and interdisciplinary professionals from the law and social science. Expertise matters, and these stakeholders can bring their varied experiences to bear on the core challenge of identifying harmful outcomes and context shifts. Technology or datasets that seem non-problematic to one group may be deemed disastrous by others. The manner in which different user groups can game certain applications or tools may also not be so obvious to the teams charged with bringing an AI-based technology to market. These kinds of impacts can sometimes be identified in early testing stages, but are usually very specific to the contextual end-use and will change over time. Acquiring these types of resources for risk and associated impacts does not necessarily require a huge allocation, but it does require deliberate planning and guidance. This is also a place where innovation in approaching bias can significantly contribute to positive outcomes.

Real-world example

There are many examples of bias from the real world where practices in the problem formulation stage may have combined with lack of understanding of downstream impacts. For example, the Gender Shades facial recognition evaluation project [24] describes the poor performance of facial recognition systems when trying to detect face types (by gender and skin type) that are not present in the training data. This is an example of representation bias – a type of sampling bias that pre-dates AI - where trends estimated for one population are inappropriately generalized to data collected from another population. This biased performance was not identified by the teams that designed and built the facial recognition systems, but instead by researchers evaluating the systems’ performance in different conditions. It is during the pre-design stage where these kinds of implicit decisions are made about what constitutes a “valid face,” and non-representative datasets are selected. Additionally, representation bias can lead to bigger problems and other biases in later stages of the AI lifecycle, an issue referred to as “error propagation,” that can eventually lead to biased outcomes [90]. Improving pre-design practices to ensure more inclusive representation can help to broaden the larger teams’ perspectives about what is considered relevant or valid.

DESIGN AND DEVELOPMENT STAGE

This stage of the AI lifecycle is where modeling, engineering and validation take place. The stakeholders in this stage tend to include software designers, engineers, and data scientists who carry out risk management techniques in the form of algorithmic auditing and enhanced metrics for validation and evaluation.

516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561

Optimization over context

The software designers and data scientists working in design and development are often highly focused on system performance and optimization. This focus can inadvertently be a source of bias in AI systems. For example, during model development and selection, modelers will almost always select the most *accurate* models. Yet, as Forde *et al* describe in their paper [50], selecting models based solely on accuracy is not necessarily the best approach for bias reduction. Not taking context into consideration during model selection can lead to biased results for sub-populations (for example, disparities in health care delivery). Relatedly, tools that are designed to use aggregated data about groups to make predictions about individual behavior – a practice initially meant to be a remedy for non-representative datasets [7]- can lead to biased outcomes. This type of bias, known as ecological fallacy, occurs when an inference is made about an individual based on their membership within a group (for example, basing college admissions decisions on an individual’s race) [48]. These unintentional weightings of certain factors can cause algorithmic results that exacerbate and reinforce societal inequities. The surfacing of these inequities is a kind of positive “side effect” of algorithmic modeling, enabling the research community to discover them and develop methods for managing them.

Practices

During modeling tasks in this stage, it may become apparent that algorithms are biased or will contribute to disparate impacts if deployed. In such cases the technology can be taken out of production. But this kind of awareness and remedy is likely to take place only in certain settings or industries, with well-defined procedures and clear lines of accountability. Unfortunately, not all tools are deployed in such settings – and capturing the wide array of use cases and scenarios is particularly difficult. It is also notable that, depending on the industry or use case, AI is typically marketed as an easy solution that does not necessarily require extensive support. But the notion that AI requires extensive monitoring belies the reality that AI can be both easy to use *and* should be used with extreme caution [96].

Several technology companies are developing or utilizing guidance to improve organizational decision making and make the practice of AI development more responsible by implementing processes such as striving to identify potential bias impacts of algorithmic models. For example, “cultural effective challenge” is a practice that seeks to create an environment where technology developers can actively challenge and question steps in modeling and engineering to help root out statistical biases and the biases inherent in human decision making [60]. Requiring AI practitioners to defend their techniques can incentivize new ways of thinking, stimulate improved practices, and help create change in approaches by individuals and organizations [96]. To better identify and mitigate organizational factors which can contribute to bias, experts also suggest the use of algorithmic decision-making tools for specific, well-defined use cases, and *not beyond* those use cases (a factor that will be discussed more in-depth in the section about deployment). Additionally, researchers also recommend that AI development teams work in tighter conjunction with subject matter experts and practitioner end users, who in turn, must “consider a deliberate and modest approach” when utilizing tool output [111].

562 Real-world example

563 One real-world case of a biased outcome that may have been manageable at the design and
564 development stage is the university admissions algorithm GRADE [135], which was shown to
565 produce biased enrollment decisions for incoming PhD students [25]. Without ground truth for
566 what constitutes a “good fit,” a construct was developed using prior admission data. Once put
567 into production, the model ended up being trained to do a different job than intended (also
568 known as “target leakage”). Instead of assessing student quality, the model learned previous
569 admissions officer decisions. Another issue is that candidate quality cannot be truly known until
570 after the student matriculates. This case is a good example of data hubris, or “overstated claims
571 that arise from big data analysis” [84]. This is particularly problematic when using data to “make
572 causal claims from an inherently inductive method of pattern recognition” [19,84,89].

573

574 DEPLOYMENT STAGE

575

576 This stage is where users start to interact with the developed technology, and sometimes create
577 unintended uses for it. The stakeholders in deployment are often the different types of end users
578 who directly interact with technology tools for their profession. This includes operators, subject
579 matter experts, humans-in-the-loop, and decision-makers who interpret output to make or
580 support decisions.

581

582 Discriminatory impact

583 Since many AI-based tools can skip deployment to a specified expert end user, and are marketed
584 to, and directly used by, the general public, the intended uses for a given tool are often quickly
585 overcome by reality. Additionally, members of the public do not necessarily have to directly
586 interact with technology to be affected by tool deployment. Individuals’ data can be used for
587 modeling (sometimes without their knowledge), and in decisions that can affect their lives based
588 on factors such as where they live and work. For example, the algorithms used in ride hailing
589 apps learned the landscape of low-income non-white neighborhoods and charged citizens who
590 live there more for pick-up and drop-off, causing disparate impact [108]. This kind of systemic
591 discriminatory pricing is perpetuated on the citizens of the neighborhood without their
592 knowledge, whether they have and use the app or not, and due only to the fact that they live
593 there.

594

595 Intended context vs. actual context

596 Once people start to interact with an AI system, early design and development decisions that
597 were poorly or incompletely specified or based on narrow perspectives can be exposed. This
598 leaves the process vulnerable to additive biases that are either statistical in nature or related to
599 human decision making and behavior [109]. For example, by not designing to compensate for
600 activity biases, algorithmic models may be built on data only from the most active users, likely
601 creating downstream system activity that does not reflect the intended or real user population
602 [1,8]. Basing system actions on an unrepresentative sample can have significant impact. For
603 example, by not considering that STEM ads might be seen most often by men, due to how
604 marketing algorithms optimize for cost in ad placement, the women who were the intended
605 audience of the ads never saw them [82].

606

607 The deployment stage also offers an interesting window into how perceptions and uses can differ
608 based on the distance from the technology itself. In pre-design the focus and perceptions are
609 about how technology can be designed to solve a question, market a product, or innovate in a
610 new area. In design and development, the focus is on building, testing, and operationalizing the
611 technology, typically with time to market and accuracy as the key criteria. And once the
612 technology is deployed and used in different settings and for different purposes, we see
613 perceptions turn to unintended use cases and even distrust. In one case of predictive analytics in
614 university admissions, the operators of the receiving end of the tool output were the ones to
615 sound the warning about race-based biases [79]. Although the study was based on a small
616 number of participants, interviews with admissions officials suggest that “they didn’t believe in
617 the validity of the risk scores, they thought the scores depersonalized their interactions with
618 students, and they didn’t understand how the scores were calculated” [48].
619

620 The kinds of scenarios where experts utilize and rely upon automated results (like in the college
621 admissions example), are highly complex and relatively understudied. One key issue is finding a
622 configuration that enables a system to be used in a way that optimally leverages, instead of
623 replaces, user expertise. This is often a significant challenge since domain experts and AI
624 developers often lack a common vernacular, which can contribute to miscommunication and
625 misunderstood capabilities. With the promise of more quantitative approaches, domain experts
626 may tend to offload method validation to the AI system itself. End users may also
627 subconsciously find ways to leverage those perceived “objective” results as cover for their biases
628 [6,38,39]. On the system side, developer communities may presume method validation at a level
629 that is not actually present. These kinds of loopholes can create conditions that operationalize
630 technology that is not quite ready for use, especially in high-stakes settings [11,120].
631

632 Contextual gaps lead to performance gaps

633 The “distance from technology” can also contribute to different types of performance gaps.
634 There are gaps in intention; these are gaps between what was originally intended in pre-design
635 versus what is developed and between the AI product and how it is deployed. There also are gaps
636 in performance based on those intention gaps. When an AI tool is designed and developed to be
637 used in a specific setting and tested for use in near-laboratory conditions, there are clearer
638 expectations about intended performance. Once the AI tool is deployed and goes “off-road,” the
639 original intent, idea, or impact assessment that was identified in pre-design can drift as the tool is
640 repurposed and/or used in unforeseen ways.
641

642 Another important gap that contributes to bias relates to differences in interpretability
643 requirements between users and developers. As previously discussed, the groups who invent and
644 produce technology have specific intentions for its use and are unlikely to be aware of all the
645 ways a given tool will be repurposed. There are individual differences in how humans interpret
646 AI model output. When system designers do not take these differences into consideration it can
647 contribute to misinterpretation of that output [21]. When these differences are combined with the
648 societal biases found in datasets and human cognitive biases such as automation complacency
649 (which is particularly relevant in the deployment stage), where end users may unintentionally
650 “offload” their decisions to the automated tool - this can cause significant negative impacts.
651
652

653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673

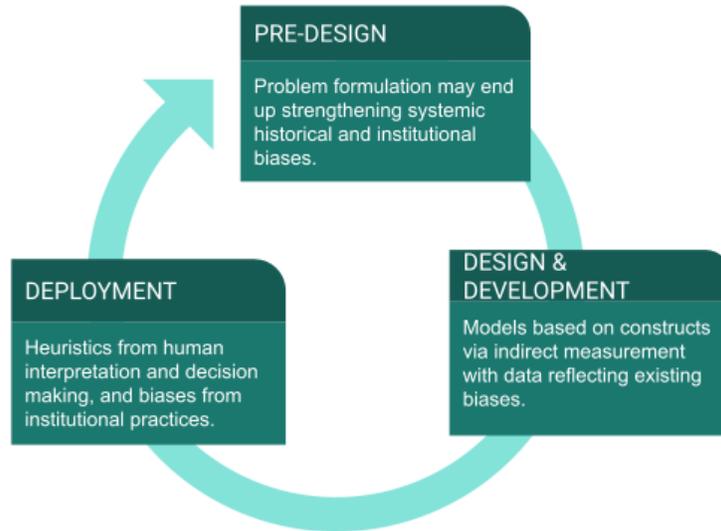
Practical improvements

One approach for managing bias risks associated with the gaps described above is deployment monitoring and auditing. Counterfactual fairness is a technique used by researchers to bridge the gaps between the laboratory and the post-deployment real world. The issue, as described in [81] is that “*If individuals in the training data have not already had equal opportunity, algorithms enforcing EO⁶ will not remedy such unfairness.*” Using the GRADE algorithm as an example, instead of using previous admission decisions as the predictor, the model would consider and seek to compensate for the various social biases that could impact a student’s application. This happens by capturing “these social biases and make clear the implicit trade-off between prediction accuracy and fairness in an unfair world.” Identifying standards of practice for implementing these types of risk management tools and techniques will be a focus of future activities.

Summary

In this section we have described the challenge of AI bias and proposed an approach for considering how to manage it through three stages modeled on AI development lifecycle. The section also shows that, while the type of bias and manner of presentation may differ, bias can occur across all of these stages. To summarize and help illustrate this point, the below figure shows an exemplar of how bias could present within each of the three stages.

Figure 2: Example of bias presentation in three stages modeled on the AI lifecycle.



674
675
676

⁶ EO = equal opportunity

677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718

5. Conclusion and Next Steps

We have identified a few of the many ways that algorithms can create conditions for discriminatory decision making. In an effort to identify the technical requirements for cultivating trustworthy and responsible AI, this report suggests a three-stage approach for managing AI bias. This approach is intended to foster discussion about the path forward and collaborative development of **standards** and a risk-based **framework**. Rather than identifying and tackling specific biases within cases, this report suggests a need to address the context-specific nature of AI bias by associating applicable biases within specific stages modeled on the AI lifecycle for more effective management and mitigation. NIST is interested in obtaining feedback from the broader community about this proposed approach via public comment and a series of public events.

The broader AI research community, practitioners, and users all have many valuable insights and recommendations to offer in managing and mitigating bias. Identifying which techniques to include in a framework that seeks to promote trustworthiness and responsibility in AI requires an approach that is actively representative and includes a broad set of disciplines and stakeholders. This will allow interested parties to move forward with guidance that is effective *and* implementable, accurate, realistic, and fit for purpose. It has the potential to increase public trust and advance the development and use of beneficial AI technologies and systems. To that end, this report concludes:

- Bias is neither new nor unique to AI.
- The goal is not zero risk but rather, identifying, understanding, measuring, managing and reducing bias.
- Standards and guides are needed for terminology, measurement, and evaluation of bias.
- Bias reduction techniques are needed that are flexible and can be applied across contexts, regardless of industry.
- NIST plans to develop a framework for trustworthy and responsible AI with the participation of a broad set of stakeholders to ensure that standards and practices reflect viewpoints not traditionally included in AI development.
- NIST will collaboratively develop additional guidance for assurance, governance, and practice improvements as well as techniques for enhancing communication among different stakeholder groups.

To make the necessary progress towards the goal of trustworthy and responsible AI, NIST intends to act as a hub for the broader community of interest and to collaboratively engage with experts and other stakeholders as they address the challenges of AI. To that end, NIST will host a variety of activities in 2021 and 2022 in each area of the core building blocks of trustworthy AI (accuracy, explainability and interpretability, privacy, reliability, robustness, safety, and security (resilience), and bias).

719 6. Appendices
 720 Appendix A: Glossary
 721 The table below presents a glossary with a stand-alone definition for each term and
 722 accompanying reference(s). The goal and contribution of this glossary is to aggregate terms that
 723 are in common usage or relevance to AI bias. Definitions were selected based on either recently
 724 published papers from the AI bias community or seminal work in the area the term is most
 725 associated with. When multiple definitions of a bias were identified, the most relevant definition
 726 was selected or adapted. The references provided are not intended to indicate specific
 727 endorsement or to assign originator credit.

728
 729 Table 1: Bias Terminology. This table lists definitions with accompanying references for select
 730 biases in AI.

Bias type	Definition
Activity bias	A type of selection bias that occurs when systems/platforms get their training data from their most active users, rather than those less active (or inactive) [8].
Amplification bias	Arises when the distribution over prediction outputs is skewed in comparison to the prior distribution of the prediction target [85].
Annotator bias, Human reporting bias	When users rely on automation as a heuristic replacement for their own information seeking and processing [93].
Automation complacency	When humans over-rely on automated systems or have their skills attenuated by such over-reliance (e.g., spelling and autocorrect or spellcheckers).
Behavioral bias	Systematic distortions in user behavior across platforms or contexts, or across users represented in different datasets [92,104].
Cognitive bias	Systematic errors in human thought based on a limited number of heuristic principles and predicting values to simpler judgmental operations [132].
Concept drift, Emergent bias	Use of a system outside the planned domain of application, and a common cause of performance gaps between laboratory settings and the real world.
Consumer bias	Arises when an algorithm or platform provides users with a new venue within which to express their biases, and may occur from either side, or party, in a digital interaction [121].
Content production bias	Arises from structural, lexical, semantic, and syntactic differences in the contents generated by users [104].
Data generation bias	Arises from the addition of synthetic or redundant data samples to a dataset [73].

Deployment bias	Arises when systems are used as decision aids for humans, since the human intermediary may act on predictions in ways that are typically not modeled in the system [127].
Detection bias	Systematic differences between groups in how outcomes are determined and may cause an over- or underestimation of the size of the effect [27].
Evaluation bias	Arises when the testing or external benchmark populations do not equally represent the various parts of the user population or from the use of performance metrics that are not appropriate for the way in which the model will be used [127].
Exclusion bias	When specific groups of user populations are excluded from testing and subsequent analyses [37].
Feedback loop bias	Effects that may occur when an algorithm learns from user behavior and feeds that behavior back into the model [121].
Funding bias	Arises when biased results are reported in order to support or satisfy the funding agency or financial supporter of the research study [91].
Historical bias	Arises when models are trained on past (potentially biased) decisions [72].
Inherited bias, Error propagation	Arises when tools that are built with machine learning are used to generate inputs for other machine learning algorithms. If the output of the tool is biased in any way, this bias may be inherited by systems using the output as input to learn other models [64].
Institutional bias, Systemic bias	A tendency for the procedures and practices of particular institutions to operate in ways which result in certain social groups being advantaged or favored and others being disadvantaged or devalued. This need not be the result of any conscious prejudice or discrimination but rather of the majority simply following existing rules or norms. Institutional racism and institutional sexism are the most common examples [28].
Interpretation bias	A form of information processing bias that can occur when users interpret algorithmic outputs according to their internalized biases and views [121].
Linking bias	Arises when network attributes obtained from user connections, activities, or interactions differ and misrepresent the true behavior of the users [104].
Loss of situational awareness bias	When automation leads to humans being unaware of their situation such that, when control of a system is given back to them in a situation where humans and machines cooperate, they are unprepared to assume their duties. This can be a loss of awareness over what automation is and isn't taking care of.
Measurement bias	Arises when features and labels are proxies for desired quantities, potentially leaving out important factors or introducing group or input-dependent noise that leads to differential performance [127].

Mode confusion bias	When modal interfaces confuse human operators, who misunderstand which mode the system is using, taking actions which are correct for a different mode but incorrect for their current situation. This is the cause of many deadly accidents, but also a source of confusion in everyday life.
Popularity bias	A form of selection bias that occurs when items that are more popular are more exposed and less popular items are under-represented [1].
Population bias	Arises when statistics, demographics, and user characteristics differ between the original target population and the user population represented in the actual dataset or platform [91].
Presentation bias	Biases arising from how information is presented on the Web, via a user interface, due to rating or ranking of output, or through users' own self-selected, biased interaction [8].
Ranking bias	The idea that top-ranked results are the most relevant and important and will result in more clicks than other results [8,86].
Sampling bias, Representation bias	Arises due to non-random sampling of subgroups, causing trends estimated for one population to not be generalizable to data collected from a new population [91].
Selection bias	Bias that results from using nonrandomly selected samples to estimate behavioral relationships as an ordinary specification bias that arises because of a missing data problem [63].
Selective adherence	Decision-makers' inclination to selectively adopt algorithmic advice when it matches their pre-existing beliefs and stereotypes [6].
Societal bias	Ascribed attributes about social groups that are largely determined by the social context in which they arise and are an adaptable byproduct of human cognition [23].
Statistical bias	A systematic tendency for estimates or measurements to be above or below their true values. Note 1: Statistical biases arise from systematic as opposed to random error. Note 2: Statistical bias can occur in the absence of prejudice, partiality, or discriminatory intent [107].
Temporal bias	Bias that arises from differences in populations and behaviors over time [104,131].
Training data bias	Biases that arise from algorithms that are trained on one type of data and do not extrapolate beyond those data.
Uncertainty bias, Epistemic uncertainty	Arises when predictive algorithms favor groups that are better represented in the training data, since there will be less uncertainty associated with those predictions [57].
User interaction bias	Arises when a user imposes their own self-selected biases and behavior during interaction with data, output, results, etc. [8].

732 Appendix B: Collaborative Work

733 This report is based on a series of collaborative events, including a literature review, input from
734 leaders in the field through ongoing discussions and a workshop, and a broad evaluation of the
735 significant themes across the community of interest. Detailed information of these events is
736 described below.

737

738 Literature review

739 During 2020, NIST implemented a broad review of materials from frequently-cited, shared, and
740 cross-referenced pieces focused on bias within technologies that use artificial intelligence. This
741 review incorporated content that described AI bias from a societal perspective, in existing
742 technologies and development processes, and other factors that influence AI development,
743 implementation, and/or adaptation. To ensure a cross-section of perspectives, literature was
744 identified across a variety of publication types, including peer-reviewed journals, popular news
745 media, books, organizational reports, conference proceedings, and presentations. Across
746 publications, the literature review topics represent a wide range of stakeholder perspectives and
747 challenges and current and future AI implementations.

748

749 Workshop on Bias in AI

750 Recognizing a lack of consensus regarding several fundamental concepts in identifying and
751 understanding bias in AI, NIST convened a virtual workshop August 18, 2020 with experts,
752 researchers, and stakeholders from a variety of organizations and sectors whose work focuses on
753 the topic. The workshop consisted of panel discussions on data and algorithmic bias, followed by
754 five contemporaneous breakout sessions. Notes from workshop organizers, facilitators, and
755 scribes were reviewed for key takeaways and themes. Workshop participants suggested that
756 forums and workshops like the one held on August 18 were important to maintaining awareness
757 and alignment of current challenges and future solutions. Participants also referred to the long-
758 term nature of this challenge. These key takeaways have been included and described throughout
759 this report.

760

761 7. References

762

- 763 [1] H. Abdollahpouri, M. Mansoury, R. Burke, B. Mobasher, The Unfairness of Popularity
764 Bias in Recommendation, ArXiv:1907.13286 [Cs]. (2019).
- 765 [2] B. Aguera y Arcas, M. Mitchell, A. Todorov, Physiognomy’s new clothes, Medium.
766 (2017).
- 767 [3] M. Aitken, E. Toreini, P. Carmichael, K. Coopamootoo, K. Elliott, A. van Moorsel,
768 Establishing a social licence for Financial Technology: Reflections on the role of the
769 private sector in pursuing ethical data practices, *Big Data & Society*. 7 (2020)
770 205395172090889.
- 771 [4] I. Ajunwa, The Paradox of Automation as Anti-Bias Intervention, *Cardozo L. Rev.* 41
772 (2020) 1671.
- 773 [5] I. Ajunwa, S. Friedler, C. Scheidegger, S. Venkatasubramanian, Hiring by Algorithm:
774 Predicting and Preventing Disparate Impact, *SSRN Electric Journal*. (2016).
- 775 [6] S. Alon-Barkat, M. Busuioc, Decision-makers Processing of AI Algorithmic Advice:
776 Automation Bias versus Selective Adherence, ArXiv:2103.02381 [Cs]. (2021).
- 777 [7] J. Angwin, J. Larson, S. Mattu, L. Kirchner, ProPublica, Machine Bias: There’s software
778 used across the country to predict future criminals. And it’s biased against blacks.,
779 ProPublica. (2016).
- 780 [8] R. Baeza-Yates, Bias on the web, *Commun. ACM*. 61 (2018) 54–61.
- 781 [9] K. Bailey, Put Away Your Machine Learning Hammer, *Criminality Is Not A Nail*, Wired.
782 (2016).
- 783 [10] J. Bajorek, Voice Recognition Still Has Significant Race and Gender Biases, *Harvard
784 Business Review*. (2019).
- 785 [11] S. Barocas, A. Biega, B. Fish, J. Niklas, L. Stark, When not to design, build, or deploy, in:
786 Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency,
787 Association for Computing Machinery, New York, NY, USA, 2020: p. 695.
- 788 [12] S. Barocas, A. Selbst, Big Data’s Disparate Impact, *California Law Review*. 104 (2016)
789 671–732.
- 790 [13] R. Bartlett, A. Morse, R. Stanton, N. Wallace, Consumer-Lending Discrimination in the
791 FinTech Era, National Bureau of Economic Research, 2019.
- 792 [14] E. Bary, How artificial intelligence could replace credit scores and reshape how we get
793 loans, *Market Watch*. (2018).
- 794 [15] R. Benjamin, *Race after technology: Abolitionist tools for the new jim code*, John Wiley
795 & Sons, 2019.
- 796 [16] M. Bogen, All the Ways Hiring Algorithms Can Introduce Bias, *Harvard Business
797 Review*. (2019).
- 798 [17] M. Bogen, A. Rieke, Help Wanted: An Examination of Hiring Algorithms, Equity, and
799 Bias, *Upturn*. (2019).
- 800 [18] M. Boyarskaya, A. Olteanu, K. Crawford, Overcoming Failures of Imagination in AI
801 Infused System Development and Deployment, ArXiv:2011.13416 [Cs]. (2020).
- 802 [19] D. Boyd, K. Crawford, CRITICAL QUESTIONS FOR BIG DATA: Provocations for a
803 cultural, technological, and scholarly phenomenon, *Information, Communication &
804 Society*. 15 (2012) 662–679.
- 805 [20] S. Brayne, Enter the Dragnet, *Logic Magazine*. (2020).

- 806 [21] D. Broniatowski, *Psychological Foundations of Explainability and Interpretability in*
807 *Artificial Intelligence*, NIST, 2021.
- 808 [22] M. Broussard, *Artificial Unintelligence: How Computers Misunderstand the World*, MIT
809 Press, 2018.
- 810 [23] R. Brown, S. Gaertner, *Blackwell Handbook of Social Psychology: Intergroup Processes*,
811 John Wiley & Sons, 2008.
- 812 [24] J. Buolamwini, T. Gebru, *Gender Shades: Intersectional Accuracy Disparities in*
813 *Commercial Gender Classification*, in: *Proceedings of Machine Learning Research*, 2018:
814 pp. 77–91.
- 815 [25] L. Burke, U of Texas will stop using controversial algorithm to evaluate Ph.D. applicants |
816 *Inside Higher Ed*, *Inside Higher Ed*. (2020).
- 817 [26] A. Caliskan, J. Bryson, A. Narayanan, *Semantics derived automatically from language*
818 *corpora contain human-like biases*, *Science*. 356 (2017) 183–186.
- 819 [27] Centre for Evidence-Based Medicine, *Catalogue of Bias*, *Catalog of Bias*. (2017).
- 820 [28] D. Chandler, R. Munday, *A Dictionary of Media and Communication*, Oxford University
821 Press, 2011.
- 822 [29] A. Chouldechova, *Fair Prediction with Disparate Impact: A Study of Bias in Recidivism*
823 *Prediction Instruments*, *Big Data*. 5 (2017) 153–163.
- 824 [30] Coalition for Critical Technology, *Abolish the #TechToPrisonPipeline*, (2020).
- 825 [31] S. Costanza-Chock, *Design Justice, A.I., and Escape from the Matrix of Domination*,
826 *Journal of Design and Science*. (2018) 96c8d426.
- 827 [32] K. Crawford, *Artificial Intelligence’s White Guy Problem*, *The New York Times*. (2016).
- 828 [33] K. Crawford, *Time to regulate AI that interprets human emotions*, *Nature*. 592 (2021)
829 167–167.
- 830 [34] C. Criado Perez, *Invisible Women: Data Bias in a World Designed for Men*, Abrams
831 Press, 2019.
- 832 [35] D. Danks, A. London, *Algorithmic Bias in Autonomous Systems*, in: *Proceedings of the*
833 *Twenty-Sixth International Joint Conference on Artificial Intelligence*, International Joint
834 *Conferences on Artificial Intelligence Organization*, Melbourne, Australia, 2017: pp.
835 4691–4697.
- 836 [36] J. Dastin, *Amazon scraps secret AI recruiting tool that showed bias against women*,
837 *Reuters*. (2018).
- 838 [37] M. Delgado-Rodriguez, *Bias*, *Journal of Epidemiology & Community Health*. 58 (2004)
839 635–641.
- 840 [38] B.J. Dietvorst, J.P. Simmons, C. Massey, *Algorithm aversion: people erroneously avoid*
841 *algorithms after seeing them err*, *J Exp Psychol Gen*. 144 (2015) 114–126.
- 842 [39] B.J. Dietvorst, J.P. Simmons, C. Massey, *Overcoming Algorithm Aversion: People Will*
843 *Use Imperfect Algorithms If They Can (Even Slightly) Modify Them*, *Management*
844 *Science*. 64 (2018) 1155–1170.
- 845 [40] C. D’Ignazio, L. Klein, *Data Feminism*, MIT Press, 2020.
- 846 [41] L. Dormehl, *Algorithms Are Great and All, But They Can Also Ruin Lives*, *Wired*.
847 (2014).
- 848 [42] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, *Fairness Through Awareness*,
849 *ArXiv:1104.3913 [Cs]*. (2011).
- 850 [43] M.C. Elish, S. Barocas, A. Plasek, K. Ferryman, *The social & economic implications of*
851 *artificial intelligence technologies in the near-term*, *AI Now*, New York, 2016.

- 852 [44] EPIC, Algorithms in the Criminal Justice System: Risk Assessment Tools, Electronic
853 Privacy Information Center, 2020.
- 854 [45] V. Eubanks, Automating inequality: How high-tech tools profile, police, and punish the
855 poor, St. Martin's Press, 2018.
- 856 [46] M. Evans, A.W. Mathews, New York Regulator Probes UnitedHealth Algorithm for
857 Racial Bias, Wall Street Journal. (2019).
- 858 [47] E. Fast, E. Horvitz, Long-Term Trends in the Public Perception of Artificial Intelligence,
859 in: Proceedings of the AAAI Conference on Artificial Intelligence, Association for the
860 Advancement of Artificial Intelligence, 2017: p. 7.
- 861 [48] T. Feathers, Major Universities Are Using Race as a "High Impact Predictor" of Student
862 Success – The Markup, The Markup. (2021).
- 863 [49] B. Fish, L. Stark, Reflexive Design for Fairness and Other Human Values in Formal
864 Models, ArXiv:2010.05084 [Cs]. (2020).
- 865 [50] J.Z. Forde, A.F. Cooper, K. Kwegyir-Aggrey, C. De Sa, M. Littman, Model Selection's
866 Disparate Impact in Real-World Deep Learning Applications, ArXiv:2104.00606 [Cs].
867 (2021).
- 868 [51] G. Friedman, T. McCarthy, Employment Law Red Flags in the Use of Artificial
869 Intelligence in Hiring, American Bar Association. (2019).
- 870 [52] H. Fry, Hello world: being human in the age of algorithms, WW Norton & Company,
871 2018.
- 872 [53] S. Furman, J. Haney, Is My Home Smart or Just Connected?, in: International Conference
873 on Human-Computer Interaction, Springer, Cham, 2020: pp. 273–287.
- 874 [54] D. Gaffney, N. Matias, Caveat emptor, computational social science: Large-scale missing
875 data in a widely-published Reddit corpus, PLOS ONE. 13 (2018) e0200162.
- 876 [55] M. Gianfrancesco, S. Tamang, J. Yazdany, G. Schmajuk, Potential Biases in Machine
877 Learning Algorithms Using Electronic Health Record Data, JAMA Intern Med. 178
878 (2018) 1544–1547.
- 879 [56] S. Goel, R. Shroff, J.L. Skeem, C. Slobogin, The Accuracy, Equity, and Jurisprudence of
880 Criminal Risk Assessment, SSRN Journal. (2018).
- 881 [57] B. Goodman, S. Flaxman, European Union regulations on algorithmic decision-making
882 and a "right to explanation," AIMag. 38 (2017) 50–57.
- 883 [58] P. Grother, M. Ngan, K. Hanaoka, Face recognition vendor test part 3: demographic
884 effects, National Institute of Standards and Technology, Gaithersburg, MD, 2019.
- 885 [59] E. Guo, K. Hao, This is the Stanford vaccine algorithm that left out frontline doctors, MIT
886 Technology Review. (2020).
- 887 [60] P. Hall, N. Gill, B. Cox, Responsible machine learning: Actionable strategies for
888 mitigating risks and driving adoption, O'Reilly Media Inc., Sebastopol, CA, 2020.
- 889 [61] M. Hardt, E. Price, N. Srebro, Equality of Opportunity in Supervised Learning,
890 ArXiv:1610.02413 [Cs]. (2016).
- 891 [62] E. Harlen, O. Schnuck, Objective or Biased, Bayerischer Rundfunk. (2021).
- 892 [63] J. Heckman, Sample Selection Bias as a Specification Error, Econometrica. 47 (1979)
893 153–161.
- 894 [64] T. Hellström, V. Dignum, S. Bensch, Bias in Machine Learning -- What is it Good for?,
895 ArXiv:2004.00686 [Cs]. (2020).
- 896 [65] M. Henry-Nickie, How artificial intelligence affects financial consumers, Brookings.
897 (2019).

- 898 [66] K. Hill, Flawed Facial Recognition Leads To Arrest and Jail for New Jersey Man, The
899 New York Times. (2020).
- 900 [67] ISO, Statistics — Vocabulary and symbols — Part 1: General statistical terms and terms
901 used in probability, 2006.
- 902 [68] ISO/IEC, Information technology — Vocabulary, International Organization for
903 Standardization, Geneva, Switzerland, 2015.
- 904 [69] ISO/IEC, Information technology — Big data — Overview and vocabulary, International
905 Organization for Standardization, Geneva, Switzerland, 2019.
- 906 [70] IT Modernization CoE, CoE Guide to AI Ethics, General Services Administration, n.d.
- 907 [71] A. Jacobs, S.L. Blodgett, S. Barocas, H. Daumé, H. Wallach, The meaning and
908 measurement of bias: lessons from natural language processing, in: Proceedings of the
909 2020 Conference on Fairness, Accountability, and Transparency, Association for
910 Computing Machinery, New York, NY, USA, 2020: p. 706.
- 911 [72] A. Jacobs, H. Wallach, Measurement and Fairness, Proceedings of the 2021 ACM
912 Conference on Fairness, Accountability, and Transparency. (2021) 375–385.
- 913 [73] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, S.-L. Kim, Communication-Efficient On-
914 Device Machine Learning: Federated Distillation and Augmentation under Non-IID
915 Private Data, ArXiv:1811.11479 [Cs, Stat]. (2018).
- 916 [74] J.E. Johndrow, K. Lum, An algorithm for removing sensitive information: Application to
917 race-independent recidivism prediction, *Ann. Appl. Stat.* 13 (2019) 189–220.
- 918 [75] F. Kamiran, A. Karim, S. Verwer, H. Goudriaan, Classifying Socially Sensitive Data
919 Without Discrimination: An Analysis of a Crime Suspect Dataset, in: 2012 IEEE 12th
920 International Conference on Data Mining Workshops, IEEE, Brussels, Belgium, 2012: pp.
921 370–377.
- 922 [76] A. Kerr, M. Barry, J. Kelleher, Expectations of artificial intelligence and the
923 performativity of ethics: Implications for communication governance, *Big Data & Society.*
924 7 (2020) 2053951720915939.
- 925 [77] L. Kirchner, M. Goldstein, Access Denied: Faulty Automated Background Checks Freeze
926 Out Renters – The Markup, The Markup. (2020).
- 927 [78] J. Kleinberg, H. Lakkaraju, J. Leskovec, J. Ludwig, S. Mullainathan, Human decisions and
928 machine predictions, *The Quarterly Journal of Economics.* 133 (2018) 237–293.
- 929 [79] S. Klempin, M. Grant, M. Ramos, Practitioner Perspectives on the Use of Predictive
930 Analytics in Targeted Advising for College Students, Community College Research
931 Center, 2018.
- 932 [80] B. Knowles, J. Richards, The Sanction of Authority: Promoting Public Trust in AI, in:
933 Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency,
934 Association for Computing Machinery, New York, NY, USA, 2021: pp. 262–271.
- 935 [81] M.J. Kusner, J.R. Loftus, C. Russell, R. Silva, Counterfactual Fairness, ArXiv:1703.06856
936 [Cs, Stat]. (2018).
- 937 [82] A. Lambrecht, C. Tucker, Algorithmic Bias? An Empirical Study into Apparent Gender-
938 Based Discrimination in the Display of STEM Career Ads, *SSRN Journal.* 65 (2019)
939 2966–2981.
- 940 [83] H. Ledford, Millions of black people affected by racial bias in health-care algorithms,
941 *Nature.* 574 (2019) 608–609.
- 942 [84] E. Lee, A. Yee, Toward Data Sense-Making in Digital Health Communication Research:
943 Why Theory Matters in the Age of Big Data, *Front. Commun.* 5 (2020) 1–11.

- 944 [85] K. Leino, E. Black, M. Fredrikson, S. Sen, A. Datta, Feature-Wise Bias Amplification,
945 ArXiv:1812.08999 [Cs, Stat]. (2019).
- 946 [86] K. Lerman, T. Hogg, Leveraging Position Bias to Improve Peer Recommendation, PLOS
947 ONE. 9 (2014) e98914.
- 948 [87] A. Liptak, Sent to Prison by a Software Program’s Secret Algorithms, The New York
949 Times. (2017).
- 950 [88] T. Maddox, J. Rumsfeld, P. Payne, Questions for Artificial Intelligence in Health Care,
951 JAMA. 321 (2019) 31–32.
- 952 [89] M. Malik, A Hierarchy of Limitations in Machine Learning, ArXiv:2002.05193 [Cs, Econ,
953 Math, Stat]. (2020).
- 954 [90] G. McGraw, H. Figueroa, V. Shepardson, R. Bonett, An architectural risk analysis of
955 machine learning systems: Toward more secure machine learning, Berryville Institute of
956 Machine Learning, Clarke County, VA, 2020.
- 957 [91] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, A. Galstyan, A Survey on Bias and
958 Fairness in Machine Learning, ArXiv:1908.09635 [Cs]. (2019).
- 959 [92] H. Miller, J. Thebault-Spieker, S. Chang, I. Johnson, L. Terveen, B. Hecht, “blissfully
960 happy” or “ready to fight”: Varying interpretations of emoji, in: Proceedings of the 10th
961 International Conference on Web and Social Media, ICWSM 2016, AAAI press, 2016: pp.
962 259–268.
- 963 [93] I. Misra, C.L. Zitnick, M. Mitchell, R. Girshick, Seeing through the Human Reporting
964 Bias: Visual Classifiers from Noisy Human-Centric Labels, in: 2016 IEEE Conference on
965 Computer Vision and Pattern Recognition (CVPR), IEEE, Las Vegas, NV, USA, 2016:
966 pp. 2930–2939.
- 967 [94] M. Mitchell, Artificial Intelligence: A Guide for Thinking Human, Farrar, Straus, and
968 Giroux, 2019.
- 969 [95] S. Mitchell, J. Shadlen, Mirror mirror: Reflections on quantitative fairness. Shira Mitchell:
970 Statistician, Shira Mitchell: Statistician. (2020).
- 971 [96] E. Moss, J. Metcalf, Ethics Owners, Data & Society. (2020).
- 972 [97] E. Moss, J. Metcalf, High Tech, High Risk: Tech Ethics Lessons for the COVID-19
973 Pandemic Response, Patterns. 1 (2020) 100102.
- 974 [98] D.K. Mulligan, J.A. Kroll, N. Kohli, R.Y. Wong, This Thing Called Fairness: Disciplinary
975 Confusion Realizing a Value in Technology, Proc. ACM Hum.-Comput. Interact. 3 (2019)
976 1–36.
- 977 [99] NIST, Framework for Improving Critical Infrastructure Cybersecurity, Version 1.1,
978 National Institute of Standards and Technology, Gaithersburg, MD, 2018.
- 979 [100] NIST, U.S. Leadership in AI: A Plan for Federal Engagement in Developing Technical
980 Standards and Related Tools, National Institute of Standards and Technology, 2019.
- 981 [101] NIST, NIST Privacy Framework: A Tool for Improving Privacy Through Enterprise Risk
982 Management, Version 1.0, National Institute of Standards and Technology, Gaithersburg,
983 MD, 2020.
- 984 [102] S.U. Noble, Algorithms of oppression: How search engines reinforce racism, NYU Press,
985 2018.
- 986 [103] Z. Obermeyer, B. Powers, C. Vogeli, S. Mullainathan, Dissecting racial bias in an
987 algorithm used to manage the health of populations, Science. 366 (2019) 447–453.
- 988 [104] A. Olteanu, C. Castillo, F. Diaz, E. Kıcıman, Social Data: Biases, Methodological Pitfalls,
989 and Ethical Boundaries, Front. Big Data. 2 (2019) 13.

- 990 [105] C. O’Neil, *Weapons of Math Destruction: How Big Data Increases Inequality and*
991 *Threatens Democracy*, Broadway Books, 2017.
- 992 [106] Organisation for Economic Co-operation and Development, *Recommendation of the*
993 *Council on Artificial Intelligence*, OECD Legal Instruments. (2019).
- 994 [107] Organization of Scientific Area Committees for Forensic Science, *OSAC Preferred Terms*,
995 *OSAC Preferred Terms*. (2021).
- 996 [108] A. Pandey, A. Caliskan, *Iterative Effect-Size Bias in Ridehailing: Measuring Social Bias*
997 *in Dynamic Pricing of 100 Million Rides*, ArXiv:2006.04599 [Cs]. (2020).
- 998 [109] S. Passi, S. Barocas, *Problem Formulation and Fairness*, *Proceedings of the Conference on*
999 *Fairness, Accountability, and Transparency*. (2019) 39–48.
- 000 [110] J. Pfeffer, K. Mayer, F. Morstatter, *Tampering with Twitter’s Sample API*, *EPJ Data Sci.* 7
001 (2018) 50.
- 002 [111] S. Picard, M. Watkins, M. Rempal, A. Kerodal, *Beyond the Algorithm: Pretrial Reform,*
003 *Risk Assessment, and Racial Fairness*, Center for Court Innovation, 2020.
- 004 [112] A. Picchi, *Job hunters face a new hurdle: Impressing AI*, CBS News. (2020).
- 005 [113] B. Plank, *What to do about non-standard (or non-canonical) language in NLP*,
006 ArXiv:1608.07836 [Cs]. (2016).
- 007 [114] C. Prunkl, C. Ashurst, M. Anderljung, H. Webb, J. Leike, A. Dafoe, *Institutionalizing*
008 *ethics in AI through broader impact requirements*, *Nature Machine Intelligence.* 3 (2021)
009 104–110.
- 010 [115] M. Raghavan, S. Barocas, *Challenges for mitigating bias in algorithmic hiring*, *Brookings.*
011 (2019).
- 012 [116] J. Redden, *The Harm That Data Do*, *Scientific American.* (2018).
- 013 [117] M. Roberts, D. Driggs, M. Thorpe, J. Gilbey, M. Yeung, S. Ursprung, A. Aviles-Rivero,
014 C. Etmann, C. McCague, L. Beer, J. Weir-McCall, Z. Teng, E. Gkrania-Klotsas, J. Rudd,
015 E. Sala, C.-B. Schönlieb, *Common pitfalls and recommendations for using machine*
016 *learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans*,
017 *Nature Machine Intelligence.* 3 (2021) 199–217.
- 018 [118] J. Sanchez-Monedero, L. Dencik, L. Edwards, *What does it mean to solve the problem of*
019 *discrimination in hiring? Social, technical and legal perspectives from the UK on*
020 *automated hiring systems*, ArXiv:1910.06144 [Cs]. (2020).
- 021 [119] H. Schellmann, *Auditors are testing hiring algorithms for bias, but big questions remain*,
022 *MIT Technology Review.* (2021).
- 023 [120] A.D. Selbst, D. Boyd, S.A. Friedler, S. Venkatasubramanian, J. Vertesi, *Fairness and*
024 *Abstraction in Sociotechnical Systems*, in: *Proceedings of the Conference on Fairness,*
025 *Accountability, and Transparency - FAT* ’19*, ACM Press, Atlanta, GA, USA, 2019: pp.
026 59–68.
- 027 [121] S. Silva, M. Kenney, *Algorithms, platforms, and ethnic bias*, *Commun. ACM.* 62 (2019)
028 37–39.
- 029 [122] T. Simonite, *How an Algorithm Blocked Kidney Transplants to Black Patients*, *Wired.*
030 (2020).
- 031 [123] M. Singh, K. Ramamurthy, *Understanding racial bias in health using the Medical*
032 *Expenditure Panel Survey data*, ArXiv:1911.01509 [Cs, Stat]. (2019).
- 033 [124] J. Sipior, *Considerations for development and use of AI in response to COVID-19*,
034 *International Journal of Information Management.* 55 (2020) 102170.
- 035 [125] A. Smith, M. Anderson, *Automation in Everyday Life*, Pew Research Center, 2017.

036 [126] M. Specia, Siri and Alexa Reinforce Gender Bias, U.N. Finds, The New York Times.
037 (2019).

038 [127] H. Suresh, J. Guttag, A Framework for Understanding Unintended Consequences of
039 Machine Learning, ArXiv:1901.10002 [Cs, Stat]. (2020).

040 [128] Y.C. Tan, L.E. Celis, Assessing Social and Intersectional Biases in Contextualized Word
041 Representations, ArXiv:1911.01485 [Cs, Stat]. (2019).

042 [129] M. Tobin, L. Matsakis, China is home to a growing market for dubious “emotion
043 recognition,” Rest of World. (2021).

044 [130] R. Tromble, Where Have All the Data Gone? A Critical Reflection on Academic Digital
045 Research in the Post-API Age, *Social Media + Society*. 7 (2021) 2056305121988929.

046 [131] Z. Tufekci, Big questions for social media big data: Representativeness, validity and other
047 methodological pitfalls, ArXiv:1403.7400 [Cs.SI]. (2014).

048 [132] A. Tversky, D. Kahneman, Judgment under Uncertainty: Heuristics and Biases, *Science*.
049 185 (1974) 1124–1131.

050 [133] W. Ware, Records, Computers and the Rights of Citizens, RAND Corporation, Santa
051 Monica, CA, 1973.

052 [134] A.L. Washington, R. Kuo, Whose side are ethics codes on?: power, responsibility and the
053 social good, in: Proceedings of the 2020 Conference on Fairness, Accountability, and
054 Transparency, ACM, Barcelona Spain, 2020: pp. 230–240.

055 [135] A. Waters, R. Miikkulainen, GRADE: Machine Learning Support for Graduate
056 Admissions, *AIMag*. 35 (2014) 64–64.

057 [136] M. Weber, M. Yurochkin, S. Botros, V. Markov, Black Loans Matter: Distributionally
058 Robust Fairness for Fighting Subgroup Discrimination, ArXiv:2012.01193 [Cs]. (2020).

059 [137] D. West, Brookings survey finds worries over AI impact on jobs and personal privacy,
060 concern U.S. will fall behind China, Brookings. (2018).

061 [138] D. West, Brookings survey finds divided views on artificial intelligence for warfare, but
062 support rises if adversaries are developing it, Brookings. (2018).

063 [139] D. West, J. Allen, Turning Point, Brookings. (2020).

064 [140] R. Wexler, When a Computer Program Keeps You in Jail, The New York Times. (2017).

065 [141] L. Wynants, B.V. Calster, G. Collins, R. Riley, G. Heinze, E. Schuit, M. Bonten, D.
066 Dahly, J. Damen, T. Debray, V. Jong, M. Vos, P. Dhiman, M. Haller, M. Harhay, L.
067 Henckaerts, P. Heus, M. Kammer, N. Kreuzberger, A. Lohmann, K. Luijken, J. Ma, G.
068 Martin, D. McLernon, C. Navarro, J. Reitsma, J. Sergeant, C. Shi, N. Skoetz, L. Smits, K.
069 Snell, M. Sperrin, R. Spijker, E. Steyerberg, T. Takada, I. Tzoulaki, S. Kuijk, B. Bussel, I.
070 Horst, F. Royen, J. Verbakel, C. Wallisch, J. Wilkinson, R. Wolff, L. Hooft, K. Moons, M.
071 Smeden, Prediction models for diagnosis and prognosis of covid-19: systematic review
072 and critical appraisal, *BMJ*. 369 (2020) m1328.

073 [142] State v. Loomis, 2016.

074