
Position: AI Agents & Liability – Mapping Insights from ML and HCI Research to Policy

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 AI agents are loosely defined as systems capable of executing complex, open-
2 ended tasks. Many have raised concerns that these systems will present significant
3 challenges to regulatory/legal frameworks, particularly in tort liability. However, as
4 there is no universally accepted definition of an AI agent, concrete analyses of these
5 challenges are limited, especially as AI systems continue to grow in capabilities.
6 In this paper, we argue that by focusing on *properties* of AI agents rather than the
7 threshold at which an AI system becomes an agent, we can map existing technical
8 research to explicit categories of “foreseeable harms” in tort liability, as well as
9 point to “reasonable actions” that developers can take to mitigate harms.

10 1 Introduction

11 AI agents are loosely defined in literature as AI systems capable of independently pursuing complex
12 goals. Existing systems, like AutoGPT [48, 4], are being enhanced with more autonomy, and future
13 agents will likely plan farther ahead, adapt and act in more complex settings. As AI agents promise
14 to support humans across a wide range of tasks, the associated reduction in human control/oversight
15 introduces notable risks and uncertainties for our legal system [18, 29]. While there is increasing
16 interest in understanding how legal and regulatory frameworks apply to the governance of AI agents
17 [28, 34, 12, 9, 41], the lack of a universally accepted definition of ‘AI agent’ limits concrete analysis.

18 In this work, we argue that by focusing on *key properties* of current and future AI systems, rather than
19 the threshold at which an AI system becomes an agent, we can already leverage Machine Learning
20 (ML) and Human-Computer Interaction (HCI) literature to provide insights on key legal questions.
21 Here, we focus on two questions in fault-based liability¹, which imposes a general duty of care to
22 avoid intentionally or negligently causing harm to others: "who is best able to prevent harms?" and
23 "what is a reasonable duty of care?"

24 In our analysis, we study on three properties (autonomy, complexity and adaptability) that measure the
25 capabilities AI agents; we organise relevant actors for liability using the AI value-chain (consisting of
26 foundation model developers upstream, application developers midstream and end-users downstream).
27 We argue that the multitude of actors in the AI value-chain, along with the increasing capabilities of
28 agents, makes it difficult to determine which actors control the outcomes of agent actions. Addressing
29 these challenges, we draw on ML and HCI research to identify harms that can arise due to the

¹Some jurisdictions have passed or proposed specific legislation on liability for AI harms. For example, the recently updated EU Product Liability Directive (European Parliament, 2024) now explicitly includes software/AI systems; EU’s proposed AI Liability Directive and California’s Bill SB1047 all target AI systems. However, traditional theories of liability still apply in most contexts.

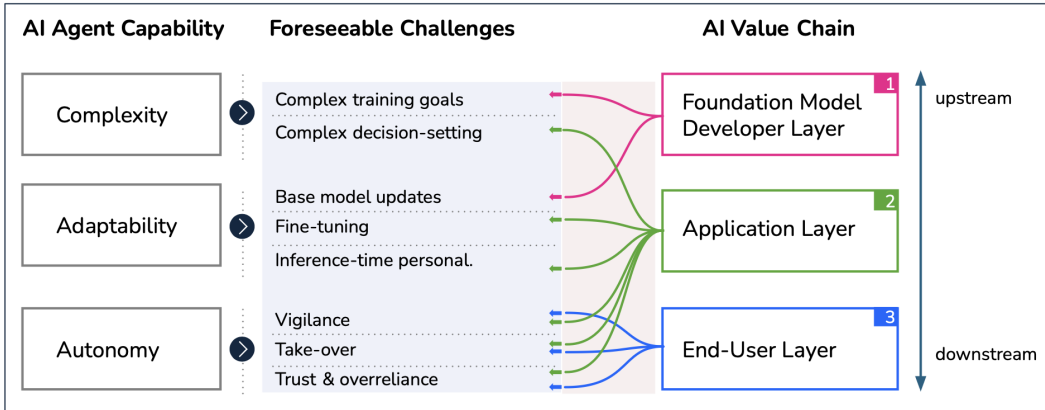


Figure 1: We map actors in the AI value chain to foreseeable challenges. Increasing agent capabilities leads to challenges (e.g. autonomy can lead to users overrelying on AI), and different actors in the chain have different degrees of control over challenges (e.g. HCI research implies that over-reliance is foreseeable by application developers). This can help define duties of care for different actors.

30 autonomy, complexity and adaptability of agents; we then identify actors who can mitigate these
 31 harms (Figure 1). Together, these insights can help define duties of care along the AI value-chain ².

32 1.1 Key Definitions

33 **What is the AI Value Chain?** We define the AI Value Chain as the sequence of steps leading to the
 34 deployment of an AI system in a specific setting. We organise the steps, as well as the actors, into
 35 layers that describe the changes made to the system. Consider the example of an AI agent deployed in
 36 a hospital that summarises patient health records for clinicians. At the Foundational Model Developer
 37 Layer, actors are developers who train and define the general capabilities of foundation models. In
 38 our example, the agent may be built upon Meta’s Meditron, a suite of open-source medical Large
 39 Language Models [10]. At the Application Layer, actors are developers who build infrastructures that
 40 allow AI systems to interface with the end-user and the environment. In our example, a third-party
 41 company may adapt the base-model, they may also build a system for the modified model to interface
 42 with hospital workers. At the End-User Layer, actors delegate tasks to an agent or use information
 43 from an agent to accomplish a specific goal. In our example, end-users are clinicians who make
 44 decisions for patient care informed by the summaries produced by the AI agent.

45 In this paper, we focus on analysing which actors in the AI Value Chain can foresee and mitigate
 46 harms that arise when the final system is deployed. In our example, consider when the AI agent
 47 produces an inaccurate summary that leads to a patient receiving an incorrect treatment and thus
 48 suffering a negative health outcome. Could the developers upstream and the clinicians downstream
 49 have foreseen this harm? What could each set of actors have done to prevent it?

50 **What is liability law?** Broadly, legal liability can arise from many areas of law, including contract
 51 law, criminal law, consumer protection law and tort law. We focus on tort liability³, which, in part,
 52 aims to ensure that victims of harmful actions are compensated, and that those responsible for causing
 53 harm are held responsible. Thus, tort liability is about corrective justice (ensuring that victims are
 54 adequately remedied), and about deterrence/harm prevention (motivating people to be careful because
 55 they can be held liable for harms they cause [46]). Based on such theories, tort liability should be
 56 placed with the person who, if acting reasonably, is able to prevent the harm.⁴

²We acknowledge that different jurisdictions have distinct tort law traditions, statutes, and case law. This paper does not aim to discuss specific legal solutions for particular jurisdictions, but instead aims for a high-level discussion on liability that may be relevant in multiple jurisdictions. We also recognize the existing legal scholarship on tort law and AI - for a comprehensive overview see footnote 5 of [25]

³Mentions of ‘liability’ in this paper should thus be read as ‘tort liability’.

⁴From a policy perspective there can also be other motivations for placing tort liability with a certain actor. For example, an actor with ‘deep pockets’ who can ensure compensation.

57 Both common law and civil law systems have ‘theories of liability’, developed through a combination
58 of jurisprudence (decisions by judges) and statutes (laws). Some common theories of liability include
59 fault liability (which includes intentional torts or negligent torts), strict liability (liability not requiring
60 ‘fault’, usually for dangerous activities or goods), product liability (liability of manufacturers for
61 defective products), and vicarious liability (liability for conduct of others). We focus on fault liability,
62 in particular on negligence, although the lessons from this paper may also be relevant for other
63 forms of tort liability or legal liability. Fault-based liability is the most general form of liability; it
64 applies even in the absence of specific legislation. For the purposes of this paper, we are focusing on
65 unintentional harms caused by AI agents, i.e. *negligent torts*.⁵

66 **What is negligence?** A ‘tortfeasor’ (person responsible for the harm) can be held liable for a negligent
67 tort when they fail to take reasonable action to prevent a foreseeable harm. Negligence hinges on the
68 concept of a breached ‘duty of care’, meaning that the potential tortfeasor did not take the reasonable
69 care expected of them in a certain situation (e.g. failing to set up a warning sign when leaving open a
70 hatch in the middle of a pedestrian walkway), and this failure leads to harm.

71 The test is not whether that specific tortfeasor had foreseen the harm (that would absolve oblivious
72 but culpable tortfeasors), but whether a ‘reasonable person’ in the position of the tortfeasor would
73 have foreseen that this harm could happen. This ‘reasonableness standard’ translates into a ‘duty of
74 care’. The duty of care is an objective standard informed by the actions of others: industry standards
75 and best practices, academic research, statements by policymakers, and legal requirements [11].

76 **What’s Challenging in Applying Existing Tort Liability Rules to AI Systems?** Value chains of AI
77 systems tend to be complex, with many actors involved in different aspects of system development
78 and deployment [5]. Under the current status quo, liability concentrates downstream towards the
79 end-user [47, 11]. This may be problematic, as downstream actors are small players and may be less
80 able to shoulder the liability compared to upstream big tech developers. Furthermore, downstream
81 actors may have less expertise, capacity or ability (e.g. access to base-models) to meet the requisite
82 duty of care [13]. Finally, AI systems may raise new risks, like immaterial harms that can result from
83 social biases baked into AI systems, that are not currently addressed by tort liability [8, 25].

84 **What is an Agent?** AI agents have been defined in different ways in literature. Generally speaking,
85 definitions center around the capability of agents: such systems would have the ability to perceive and
86 operate in complex environments, and to autonomously adapt their strategies and actions based on
87 new input [35, 21, 40, 23, 15]. In this paper, we define an AI agent as *an AI system that is deployed*
88 *in a real-life decision making setting*, and we focus on the capability of these system as measured by
89 autonomy, adaptability and complexity. We focus on these properties as they significantly challenge
90 humans’ ability to anticipate/prevent the harmful outcomes of AI agent actions.

91 **2 Diffusion of Liability Along the AI Value Chain**

92 We describe specific ways the autonomy, adaptability and complexity of AI agents may lead to
93 harmful outcomes. Furthermore, for each actor in the AI value chain, we draw from current research
94 to identify the degrees to which they may foresee and mitigate harms (summary in Figure 1).

95 **Autonomy.** We focus on two levels of autonomy when AI agents interact with human users: (1) AI
96 agents that support human decision-making, and (2) AI agents that operate under human supervision.

97 AI agents with low autonomy include current decision-support tools that rely on humans for critical
98 assessment of AI outputs. For example, AI agents have been deployed to support clinical decision
99 making, such as by performing risk assessments of patients and recommending treatment [38].
100 Clinicians using these systems need to determine whether or not the AI recommendations are valid or
101 useful. However, literature in Human-Computer Interaction (HCI) demonstrates that humans have a
102 strong propensity to over-rely on and over-trust AI assistants [22, 6].

103 For AI agents with greater degrees of autonomy, a common design choice is to cast human users
104 in supervisory roles, monitoring the system for errors and taking over during exceptional circum-
105 stances [3]. In a clinical setting, this may look like an agent that automatically screens mammography
106 for cancer, supervised by a clinical expert who monitors the system and steps in to perform manual

⁵Liability is relatively straightforward for intentional harms, as it will usually be the person intentionally causing the harm who ‘controlled the outcome’ and is liable.

107 diagnosis in cases of errors and exceptions. Unfortunately, a large body of work in HCI shows us
108 that humans struggle with vigilance (i.e. sustaining attention while performing monotonous tasks
109 over long periods of time), thus making them poor supervisors of automated systems [1, 42, 37]. In
110 exceptional cases where humans are required to take-over for the AI agent, the problem of vigilance
111 amplifies the difficulty of control transfer: the human, having reduced prior involvement with the
112 task, is likely to struggle with reacting appropriately and in a timely fashion to the situation [31].

113 *What are 'foreseeable harms'*: Overwhelmingly, literature on human factors in computing teaches us
114 that naïve integration of agents into AI+human teams can exacerbate human error [2, 50]. That is,
115 existing HCI research can help establish explicit and concrete categories of 'foreseeable harms', as
116 well as point to 'reasonable actions' application developers can take to mitigate these harms.

117 *Who can prevent these harms*: Design of human+AI interfaces happens in the Application Layer,
118 where developers can leverage literature to address problematic ways end-users will interact with
119 agents. For example, AI systems that provide continuous-support (rather than recommendation-centric
120 support) have been shown to help users maintain situational awareness in human+AI teams [52].

121 **Adaptability.** We focus on two common ways that agents adapt: (1) changing the system's base-
122 model at pre-deployment or model-update time, (2) personalising system output with specialised
123 input at inference (i.e. decision-making) time.

124 Many foundation models can be customised. Outside of open-source models, major developers like
125 OpenAI, Google, Microsoft, Meta, Anthropic, and Amazon have existing or proposed mechanisms
126 for downstream developers to adapt their models for specific applications through fine-tuning [33].
127 That is, these models can be updated based on new data and interactions in order to learn new skills.
128 However, unlike traditional software updates, which are designed to preserve existing functionality
129 whilst adding new ones, updating AI models with new data can lead to (1) degradation of existing
130 functionality [51] (for example, safety guardrails can be quickly by-passed with fine-tuning, providing
131 user access to dangerous information); (2) unexpected new biases and failure modes, as new data
132 interacts with existing ones on which the model was trained [45].

133 Even when models in AI agents are fixed, the outputs of these systems can be personalised to
134 individual or groups of end-users at inference-time, by including special information in system inputs.
135 For example, language models can suggest personalised email subjects when shown user's past
136 emails [36]. However, personalising model outputs risks confirmation bias (selective reinforcement of
137 users' existing opinions) [39, 14]; it can even result in behaviours that can be categorised as deceptive,
138 e.g. actively steering users away from or hiding contradictory information [20]. Furthermore, and
139 perhaps surprisingly, personalisation can harm model performance. That is, personalising a model to
140 a specific group of users can lower the model's performance at a group level [43, 53]. For example,
141 providing gender to language models when generating a recommendation letter can increase model
142 hallucination as well as diminish language associated to "excellence" [44].

143 *What are 'foreseeable harms'*: Research shows that, for adaptable agents, anticipating harms based on
144 pre-deployment functionalities cannot cover the range of harms that may result from post-deployment
145 changes. Existing works that study AI models in the regimes of fine-tuning, continual learning and
146 transfer learning can already help us anticipate emergent agent behaviours and associated risks.

147 *Who can prevent these harms*: How much and in what ways an agent can adapt are design choices
148 made both at the Foundation Model Developer Layer and the Application Layer. While safe and
149 effective adaptation is an open research question, upstream developers have a responsibility to test
150 agents for failures known to arise in adaptation, and disclose these risks to downstream actors.

151 **Complexity.** We focus on two aspects of complexity: (1) complexity of the agent's goals (specifically,
152 its training objective), and (2) complexity of the agent's decision setting (specifically, the length of
153 the planning horizon and the effective size of the environment).

154 Behaviours of AI agents are often determined by multiple objectives that may be in tension. For
155 example, foundation models trained with human feedback implicitly balance potentially conflicting
156 preferences of different users. Furthermore, there are tensions between social welfare goals (e.g.
157 safety) and personalisation goals (e.g. open access to information for individual users). However,
158 agents are often trained by maximising a single combination of multiple objectives, without explicitly
159 managing trade-offs amongst them. Thus, the resulting agent can make unexpected and undesirable
160 compromises [19], for example, by sacrificing safety in order to align with user preferences.

161 In addition to the multiple objectives, models in AI agents are frequently trained with “blackbox” and
162 underspecified objectives. For example, foundation models are often fine-tuned with direct human
163 feedback on model outputs [27], or with datasets of human annotated examples that encode human
164 preference [24]. The objective implied by human feedback and annotated examples is not defined
165 in closed-form, and hence cannot be directly inspected or interrogated. Furthermore, given a set of
166 examples, there are often multiple plausible objectives that a model can infer that would cause it to
167 generalize in dramatically different ways. Thus, training with implicit and underspecified objectives
168 often leads to unexpected and undesirable agent behaviours [30, 26]. For example, an agent trained
169 to interact with humans in more naturalistic and context-sensitive ways can learn, as unintended
170 side-effects, to hold strong political opinions and to pursue potentially dangerous “subgoals” (e.g. to
171 accumulate resources for the pursuit of current goal) [32].

172 Even when training objectives of AI agents can be validated, it is hard to design interfaces that
173 allow end-users to anticipate outcomes of agent actions, when actions take place over long planning
174 horizons and in open environments. Existing works show that humans can even struggle to understand
175 decisions of simple AI systems in single-shot decision settings [7]. As the planning horizons of agents
176 increase, and as the size of the environment as well as the number of other agents in the environment
177 grow, the sequence of agent decisions becomes too complex to directly inspect. For example, an open
178 question is how to effectively summarise complex policies for sequential decision-making: simply
179 enumerating the agent’s actions over the large number of possible states of the environment yields
180 results that are uninterpretable to end-users [49, 16].

181 *What are ‘foreseeable harms’:* A number of failures of AI models can be exposed by quantifying how
182 trained models trade-off different task-relevant desiderata. There is also a large body of literature that
183 can surface biases models learn from human data [17]. Finally, research on explainable AI (XAI), as
184 well as HCI, can anticipate failures of Human+AI teams in complex decision settings.

185 *Who can prevent these harms:* The complexity of an agent’s goals can be determined at the Founda-
186 tional Model Developer Layer (where the capabilities of the base-model is determined) as well as at
187 the Application Layer (where capabilities may be added). Upstream developers have a responsibility
188 to explicitly prevent models from making undesirable trade-offs. The complexity of an agent’s
189 decision setting is primarily determined at the Application Layer. Application developers need to
190 design Human+AI interfaces guided by best practices in XAI and HCI.

191 3 Implications for Policy & Technical Research

192 Increasing capabilities of AI agents challenge human control when these systems are deployed in
193 real-life. However, we argue that existing research in ML and HCI point to ways that actors in the
194 AI value chain can already better foresee and mitigate potential harms. From the perspective of
195 fault-based liability, actors at the *Foundation Model Developer Layer* have a responsibility to disclose
196 models’ training objectives and understand trade-offs models make between different objectives.
197 They have a responsibility to check for and mitigate known model biases that arise from data selection
198 and training; they should also test for how biases affect models in common down-stream tasks.
199 Safe-guards should be implemented against inappropriate and potentially unsafe types of model
200 customisation. At the *Application Layer*, when customising models, developers have a responsibility
201 to test for (and address) known model failures due to fine-tuning, as well as for unequal model
202 performance when it is personalised to different end-users. Furthermore, when integrating models
203 into Human+AI systems, they have a responsibility to anticipate/address known challenges in human-
204 AI interactions. At the *End-User Layer*, users should understand the types of tasks the AI agent can
205 safely perform. They should especially monitor the agent’s behaviour when it is used in a new setting.

206 Upstream actors have a general responsibility to expose conditions under which the model has been
207 tested, and how it scored on different evaluations. Transparency around evaluation settings helps
208 downstream actors perform due diligence when choosing to deploy AI agents for a specific use-case.

209 **Call for interdisciplinary research.** We see a need for more interdisciplinary research between core
210 ML and HCI, to concretely connect properties of models (e.g. catastrophic forgetting, interference,
211 reward hacking, pathologies arising from multi-objective optimization) to specific impacts on users’
212 abilities to anticipate/mitigate harms of AI systems. We also see a need for more collaborations
213 between policy/law makers and technical researchers to formalise regulatory/legal principles as
214 technical desiderata (e.g. training procedures, objectives and metrics, socio-technical evaluations).

215 **Impact Statement**

216 We argue that interdisciplinary collaboration between technical (especially between ML and HCI)
217 researchers and policy/law makers is necessary for effectively and responsibly addressing emerging
218 societal challenges due to complex AI systems. By involving ML and HCI researchers in policy/legal
219 processes, policy and lawmakers can better understand capabilities and limitations of emerging
220 technologies, enabling the development of regulations that are informed and effective. For the
221 technical community, engagement with law/policy professionals opens avenues for shaping regulatory
222 environments in which they operate. Collaborating with legal/policy experts ensures that ML systems
223 are designed with compliance in mind, reducing the risk of costly redesigns or legal challenges. It
224 also encourages the adoption of ethical best practices, fostering public trust in AI applications.

225 **A Appendix: Recommendations for Policymakers**

226 As AI agents grow in autonomy, complexity, and adaptability, assigning and determining liability
227 for harms becomes increasingly challenging. In this paper, we have shown that interdisciplinary
228 research, particularly at the intersection of Machine Learning (ML), Human-Computer Interaction
229 (HCI), and law, can help us establish clearer liability frameworks. Achieving this, we believe that
230 there are concrete actions that policymakers can take to create the foundation for more robust and
231 nuanced research into AI liability. We recommend the following:

232 **Enhance visibility into AI agent development and deployment**

- 233 • Facilitate research access to agent usage data to better understand human-agent interactions,
234 similar to research access provisions under the Digital Services Act (DSA)
- 235 • Require identification systems for agents [?] and logging mechanisms to verify AI agents’
236 actions
- 237 • Mandate stage- and context-specific evaluations and audits throughout the AI agent lifecycle
238 through comprehensive auditing regimes
- 239 • Establish mandatory incident reporting and information sharing protocols

240 **Strengthen technical governance capacity to support judicial assessment of liability**

- 241 • Ensure sufficient socio-technical expertise within regulatory agencies through targeted
242 hiring, training programs, and continuous education initiatives
- 243 • Explore flexible models such as scientific advisory boards or fellowship programs to source
244 and integrate external expertise
- 245 • Develop standardized methodologies for assessing AI system capabilities and limitations in
246 legal contexts

247 **Establish and enforce a comprehensive duty of care for AI development and deployment**

- 248 • Incentivize the development and adoption of industry-wide standards, particularly focusing
249 on safety measures and trust infrastructures
- 250 • Support research to concretely connect properties of AI models (e.g., catastrophic forgetting,
251 interference, reward hacking) to legal liability frameworks
- 252 • Facilitate the public dissemination of new safety procedures and formally promulgate these
253 standards through regulatory channels

254 **References**

- 255 [1] Julia Adler-Milstein, Donald A Redelmeier, and Robert M Wachter. The limits of clinician
256 vigilance as an ai safety bulwark. *Jama*, 2024.
- 257 [2] Saar Alon-Barkat and Madalina Busuioc. Human–ai interactions in public sector decision
258 making: “automation bias” and “selective adherence” to algorithmic advice. *Journal of Public
259 Administration Research and Theory*, 33(1):153–169, 2023.

- 260 [3] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamvi-
261 boonsuk, and Laura M Vardoulakis. A human-centered evaluation of a deep learning system
262 deployed in clinics for the detection of diabetic retinopathy. In *Proceedings of the 2020 CHI*
263 *conference on human factors in computing systems*, pages 1–12, 2020.
- 264 [4] Timo Birr, Christoph Pohl, Abdelrahman Younes, and Tamim Asfour. Autogpt+ p: Affordance-
265 based task planning with large language models. *arXiv preprint arXiv:2402.10778*, 2024.
- 266 [5] Ian Brown. Allocating accountability in ai supply chains, 2023.
- 267 [6] Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z Gajos. To trust or to think: cognitive
268 forcing functions can reduce overreliance on ai in ai-assisted decision-making. *Proceedings of*
269 *the ACM on Human-computer Interaction*, 5(CSCW1):1–21, 2021.
- 270 [7] Adrian Bussone, Simone Stumpf, and Dymna O’Sullivan. The role of explanations on trust and
271 reliance in clinical decision support systems. In *2015 international conference on healthcare*
272 *informatics*, pages 160–169. IEEE, 2015.
- 273 [8] Alan Chan, Rebecca Salganik, Alva Markelius, Chris Pang, Nitarshan Rajkumar, Dmitrii
274 Krasheninnikov, Lauro Langosco, Zhonghao He, Yawen Duan, Micah Carroll, et al. Harms
275 from increasingly agentic algorithmic systems. In *Proceedings of the 2023 ACM Conference on*
276 *Fairness, Accountability, and Transparency*, pages 651–666, 2023.
- 277 [9] Alan Chan, Carson Ezell, Max Kaufmann, Kevin Wei, Lewis Hammond, Herbie Bradley, Emma
278 Bluemke, Nitarshan Rajkumar, David Krueger, Noam Kolt, et al. Visibility into ai agents. In
279 *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 958–973,
280 2024.
- 281 [10] Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba,
282 Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami,
283 et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint*
284 *arXiv:2311.16079*, 2023.
- 285 [11] Peter Cihon. Chilling autonomy: Policy enforcement for human oversight of ai agents. In *41st*
286 *International Conference on Machine Learning, Workshop on Generative AI and Law*, 2024.
- 287 [12] Michael K Cohen, Noam Kolt, Yoshua Bengio, Gillian K Hadfield, and Stuart Russell. Regulat-
288 ing advanced artificial agents. *Science*, 384(6691):36–38, 2024.
- 289 [13] Connor Dunlop. Regulating ai foundation models is crucial for innovation, 2023.
- 290 [14] Sharon Ferguson, Paula Akemi Aoyagui, Young-Ho Kim, and Anastasia Kuzminykh. Just
291 like me: The role of opinions and personal experiences in the perception of explanations in
292 subjective decision-making. *arXiv preprint arXiv:2404.12558*, 2024.
- 293 [15] Iason Gabriel, Arianna Manzini, Geoff Keeling, Lisa Anne Hendricks, Verena Rieser, Hasan
294 Iqbal, Nenad Tomašev, Ira Ktena, Zachary Kenton, Mikel Rodriguez, et al. The ethics of
295 advanced ai assistants. *arXiv preprint arXiv:2404.16244*, 2024.
- 296 [16] Jasmina Gajcin, Rahul Nair, Tejaswini Pedapati, Radu Marinescu, Elizabeth Daly, and Ivana
297 Dusparic. Contrastive explanations for comparing preferences of reinforcement learning agents.
298 *arXiv preprint arXiv:2112.09462*, 2021.
- 299 [17] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck
300 Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language
301 models: A survey. *Computational Linguistics*, pages 1–79, 2024.
- 302 [18] Sara Gerke, Timo Minssen, and Glenn Cohen. Ethical and legal challenges of artificial
303 intelligence-driven healthcare. In *Artificial intelligence in healthcare*, pages 295–336. Elsevier,
304 2020.
- 305 [19] Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane,
306 Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz,
307 et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous*
308 *Agents and Multi-Agent Systems*, 36(1):26, 2022.

- 309 [20] Betty Li Hou, Kejian Shi, Jason Phang, James Aung, Steven Adler, and Rosie Campbell. Large
310 language models as misleading assistants in conversation. *arXiv preprint arXiv:2407.11789*,
311 2024.
- 312 [21] Qiuyuan Huang, Naoki Wake, Bidipta Sarkar, Zane Durante, Ran Gong, Rohan Taori, Yusuke
313 Noda, Demetri Terzopoulos, Noboru Kuno, Ade Famoti, et al. Position paper: Agent ai towards
314 a holistic intelligence. *arXiv preprint arXiv:2403.00833*, 2024.
- 315 [22] Sarah Jabbour, David Fouhey, Stephanie Shepard, Thomas S Valley, Ella A Kazerooni, Nikola
316 Banovic, Jenna Wiens, and Michael W Sjoding. Measuring the impact of ai in the diagnosis of
317 hospitalized patients: a randomized clinical vignette survey study. *Jama*, 330(23):2275–2284,
318 2023.
- 319 [23] Sayash Kapoor, Benedikt Stroebel, Zachary S Siegel, Nitya Nadgir, and Arvind Narayanan. Ai
320 agents that matter. *arXiv preprint arXiv:2407.01502*, 2024.
- 321 [24] Y Kim, X Xu, D McDuff, C Breazeal, and HW Park. Health-llm: large language models for
322 health prediction via wearable sensor data, arxiv. *arXiv preprint arXiv:2401.06866*, 2024.
- 323 [25] Noam Kolt. Governing ai agents. *Available at SSRN*, 2024.
- 324 [26] Victoria Krakovna. Specification gaming examples in AI, 2018.
- 325 [27] Hengli Li, Song-Chun Zhu, and Zilong Zheng. Diplomat: a dialogue dataset for situated
326 pragmatic reasoning. *Advances in Neural Information Processing Systems*, 36:46856–46884,
327 2023.
- 328 [28] Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca
329 Dragan. Learning to model the world with language. In Ruslan Salakhutdinov, Zico Kolter,
330 Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp,
331 editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235
332 of *Proceedings of Machine Learning Research*, pages 29992–30017. PMLR, 21–27 Jul 2024.
333 URL <https://proceedings.mlr.press/v235/lin24g.html>.
- 334 [29] John D McGreevey, C William Hanson, and Ross Koppel. Clinical, legal, and ethical aspects
335 of artificial intelligence–assisted conversational agents in health care. *Jama*, 324(6):552–553,
336 2020.
- 337 [30] Alexander Pan, Erik Jones, Meena Jagadeesan, and Jacob Steinhardt. Feedback loops with
338 language models drive in-context reward hacking. *arXiv preprint arXiv:2402.06627*, 2024.
- 339 [31] Samir Passi and Mihaela Vorvoreanu. Overreliance on ai literature review. *Microsoft Research*,
340 2022.
- 341 [32] Ethan Perez, Sam Ringer, Kamilė Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig
342 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model
343 behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022.
- 344 [33] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson.
345 Fine-tuning aligned language models compromises safety, even when users do not intend to!
346 *arXiv preprint arXiv:2310.03693*, 2023.
- 347 [34] Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton,
348 Bethanie Brownfield, Gavin Buttimore, Max Cant, Sarah Chakera, et al. Scaling instructable
349 agents across many simulated worlds. *arXiv preprint arXiv:2404.10179*, 2024.
- 350 [35] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson, 2016.
- 351 [36] Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. Lamp: When large
352 language models meet personalization. *arXiv preprint arXiv:2304.11406*, 2023.
- 353 [37] Abigail Sellen and Eric Horvitz. The rise of the ai co-pilot: Lessons for design from aviation
354 and beyond. *Communications of the ACM*, 67(7):18–23, 2024.

- 355 [38] Mark P Sendak, William Ratliff, Dina Sarro, Elizabeth Alderton, Joseph Futoma, Michael Gao,
356 Marshall Nichols, Mike Revoir, Faraz Yashar, Corinne Miller, et al. Real-world integration of a
357 sepsis deep learning technology into routine clinical care: implementation study. *JMIR medical*
358 *informatics*, 8(7):e15182, 2020.
- 359 [39] Nikhil Sharma, Q Vera Liao, and Ziang Xiao. Generative echo chamber? effect of llm-powered
360 search systems on diverse information seeking. In *Proceedings of the CHI Conference on*
361 *Human Factors in Computing Systems*, pages 1–17, 2024.
- 362 [40] Yonadav Shavit, Sandhini Agarwal, Miles Brundage, Steven Adler, Cullen O’Keefe, Rosie
363 Campbell, Teddy Lee, Pamela Mishkin, Tyna Eloundou, Alan Hickey, et al. Practices for
364 governing agentic ai systems. *Research Paper, OpenAI, December*, 2023.
- 365 [41] Julia Smakman. Ai assistants: Helpful or full of hype?, 2024.
- 366 [42] Helen Smith and Kit Fotheringham. Artificial intelligence in clinical decision-making: rethink-
367 ing liability. *Medical Law International*, 20(2):131–154, 2020.
- 368 [43] Vinith Menon Suriyakumar, Marzyeh Ghassemi, and Berk Ustun. When personalization harms
369 performance: reconsidering the use of group attributes in prediction. In *International Conference*
370 *on Machine Learning*, pages 33209–33228. PMLR, 2023.
- 371 [44] Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. "kelly
372 is a warm person, joseph is a role model": Gender biases in llm-generated reference letters.
373 *arXiv preprint arXiv:2310.09219*, 2023.
- 374 [45] Zirui Wang, Zachary C Lipton, and Yulia Tsvetkov. On negative interference in multilingual
375 models: Findings and a meta-learning treatment. *arXiv preprint arXiv:2010.03017*, 2020.
- 376 [46] Ernest J Weinrib. Deterrence and corrective justice. *UCLA L. Rev.*, 50:621, 2002.
- 377 [47] Maria Yagoda. Airline held liable for its chatbot giving passenger bad advice—what this
378 means for travelers. URL [https://www.bbc.com/travel/article/20240222-air-canada-chatbot-](https://www.bbc.com/travel/article/20240222-air-canada-chatbot-misinformation-what-travellers-should-know)
379 [misinformation-what-travellers-should-know](https://www.bbc.com/travel/article/20240222-air-canada-chatbot-misinformation-what-travellers-should-know), 2024.
- 380 [48] Hui Yang, Sifu Yue, and Yunzhong He. Auto-gpt for online decision making: Benchmarks and
381 additional opinions. *arXiv preprint arXiv:2306.02224*, 2023.
- 382 [49] Jiayu Yao, Sonali Parbhoo, Weiwei Pan, and Finale Doshi-Velez. Policy optimization with
383 sparse global contrastive explanations. *arXiv preprint arXiv:2207.06269*, 2022.
- 384 [50] Hubert D Zając, Dana Li, Xiang Dai, Jonathan F Carlsen, Finn Kensing, and Tariq O Andersen.
385 Clinician-facing ai in the wild: Taking stock of the sociotechnical challenges and opportunities
386 for hci. *ACM Transactions on Computer-Human Interaction*, 30(2):1–39, 2023.
- 387 [51] Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma.
388 Investigating the catastrophic forgetting in multimodal large language models. *arXiv preprint*
389 *arXiv:2309.10313*, 2023.
- 390 [52] Zelun Tony Zhang, Sebastian S Feger, Lucas Dullenkopf, Rulu Liao, Lukas Süßlin, Yuanting
391 Liu, and Andreas Butz. Beyond recommendations: From backward to forward ai support of
392 pilots’ decision-making process. *arXiv preprint arXiv:2406.08959*, 2024.
- 393 [53] Mingqian Zheng, Jiaxin Pei, and David Jurgens. Is "a helpful assistant" the best role for large
394 language models? a systematic evaluation of social roles in system prompts. *arXiv preprint*
395 *arXiv:2311.10054*, 2023.