### Mechanisms for *Aggregated Individual Reporting* Should Be Established for Post-Deployment Evaluation

Anonymous Authors<sup>1</sup>

#### Abstract

The need for developing model evaluations beyond static benchmarking is now well-understood, bringing attention to post-deployment auditing and evaluation. At the same time, concerns about the concentration of power in deployed AI systems have sparked a keen interest in "democratic" or "public" AI. In this work, we bring these two ideas together by proposing mechanisms for aggregated individual reporting (AIR), a framework for post-deployment evaluation that relies on individual reports from the public. AIRs allow those who interact with a specific, deployed (AI) system to report when they feel that they may have experienced something problematic; these reports are then aggregated over time, with the goal of evaluating the relevant system in a fine-grained manner. This position paper argues that individual experiences should be understood as an integral part of post-deployment evaluation, and that the scope of our proposed aggregated individual reporting mechanism is a practical path to that end. On the one hand, individual reporting can identify substantively novel insights about safety and performance; on the other, aggregation can be uniquely useful for informing action. From a normative perspective, the post-deployment phase completes a missing piece in the conversation about "democratic" AI. As a pathway to implementation, we provide a workflow of concrete design decisions and pointers to areas requiring further research and methodological development.

#### 1. Introduction

In the third week of April 2025, OpenAI quietly pushed an update to GPT-40, the model that powers the ChatGPT product by default. Online, the complaints started rolling in: The newest update was annoying and introduced bugs; it overpromised and underdelivered. More frighteningly, it encouraged users to stop taking their medications, validated conspiracy theories, and worse. By April 29, around one week later, OpenAI had rolled back the change (OpenAI, 2025b).

In this case, the unstructured feedback of individual users was essential. OpenAI was unable to identify ahead of time that this "personality" update was problematic—in part because it is hard to anticipate the richness of usage patterns, and therefore failure modes. Users thought the problems were egregious enough that it motivated them to share online; enough users tweeted about the *same* problem that OpenAI noticed. This ChatGPT personality problem was serious, widespread, and was therefore caught even in the absence of a formal system to collect feedback. But what other, subtler, non-Twitter-viral patterns might be happening—and how might we find out?

In this work, we formalize *aggregated individual reporting* (AIR) as a general framework for understanding the realworld impact of an AI system in a structured, intentional, and thorough manner.<sup>1</sup> At a high level, an AIR allows those who interact with a specific AI system to submit feedback (reports) when they believe they have experienced harm due to the system. The mechanism aggregates these reports over time, with the goal of building collective knowledge about the contours of system behavior. Intuitively, one person having one bad experience does not by itself necessarily imply a system-level problem. On the other hand, if many reports of similar experiences begin to accumulate, then perhaps there may be an important underlying issue that requires action.

Our framework relies on two crucial beliefs: first, that those

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

<sup>Preliminary work. Under review by the Workshop on Models of
Human Feedback for AI Alignment at the International Conference</sup> on Machine Learning (ICML). Do not distribute.

<sup>&</sup>lt;sup>1</sup>Throughout this paper, we will often use "AIR" as a stand-in for "aggregated individual reporting mechanism" for concision and clarity, and to distinguish the full mechanism from its components (aggregation and individual reporting).

interacting with the system have unique and valuable perspectives on its failure modes, and second, that they ought

57 to be able to express those perspectives in a nontrivial way.

This position paper argues that post-deployment evalua-

059 tion must account for individual experiences—and that

060 aggregated individual reporting mechanisms, as we define

061 in this work, are a practical pathway to doing so.

062 We do not claim to originate the concept of individual re-063 porting. In fact, we are inspired by the existence of similar 064 reporting mechanisms in other domains, like the Vaccine 065 Adverse Events Reporting System (VAERS) (Shimabukuro 066 et al., 2015), the Aviation Safety Reporting System (ASRS) 067 (Beaubien & Baker, 2002), and various medical reporting 068 systems (Wu et al., 2002). We are also heavily influenced by 069 calls to incorporate end-user expertise in algorithm evalua-070 tions (e.g., Shen et al. (2021); DeVos et al. (2022); Lam et al. (2022)) and to shift power towards the public (e.g., Kalluri 072 (2020); Feffer et al. (2023)). Moreover, while several recent 073 policy directives mandate the consideration of public feedback with specific reference to post-deployment monitoring 075 (e.g., the EU AI Act (European Parliament, 2024), the UN 076 General Assembly's first AI resolution (Assembly, 2024), 077 and Biden's now-repealed executive order (Biden, 2023)), 078 they are, by nature, only high level. By presenting a high-079 level structure for AIRs, we aim to offer a shared language as starting point for conversation across subfields and do-081 mains, and to concretize what effective implementation of 082 this policy might look like. 083

The remainder of this paper is organized as follows. In 085 Section 2, we define our idealized version of an AIR mech-086 anism, and contextualize our definition with the current 087 state-of-the-art for post-deployment evaluation and audit-088 ing for AI systems. In Section 3, we argue that individual 089 reporting effectively identifies "unknown unknowns," and 090 that aggregation enables downstream action. In Section 4, 091 we elaborate on the normative case for individual report-092 ing, placing our proposal in the context of recent calls for 093 "democratic" AI. Section A describes more granular design 094 decisions as well as associated open research questions. We 095 conclude by discussing of key challenges for successful real-world deployment in Section 5. 096

## 0980992. Defining aggregated individual reporting

We begin by describing our idealized vision of a mechanism
for aggregated individual reporting, illustrated in Figure 1.
This definition is intentionally broad, so as to cover a wide
range of potential applications—see Section A for more
granular details that describe various potential implementations. Even with this broad scope, we note that, to the best
of our knowledge, (public or non-proprietary) AIRs for AI
systems are not currently in mainstream use.

### 108

097

## 2.1. Defining an idealized individual reporting mechanism

Our understanding of an AIR mechanism begins with three key entities.

First, the *evaluated system* determines the scope of the AIR mechanism. For example, this could be a general-purpose model like ChatGPT or Claude; a product that scribes inperson doctors' appointments; or a predictive algorithm that several banks use for making loan decisions. Reports are submitted about the evaluated system. Second, the evaluated system induces an *affected population*. This could include users of ChatGPT, or their loved ones; patients and healthcare workers; or loan applicants. Reports are submitted by members of the affected population. Finally, the AIR mechanism is orchestrated by a mechanism administrator, which is the organizational entity that collects and aggregates the reports. The administrator makes key decisions like determining the scope of the evaluated system and the affected population, as well as various implementation details as discussed in Section A. The mechanism administrator could be first-party, i.e. the same organization that operates the evaluated system (e.g., OpenAI collecting reports for ChatGPT); second-party, i.e. an organizational user of the evaluated system (e.g., a hospital system collecting reports about an AI scribe product); or third-party, i.e. an external organization (e.g., a government body or activist nonprofit collecting reports about a loan allocation algorithm). This taxonomy matches prior work in the AI audit space (e.g., Raji et al. (2022)).

AIR mechanisms comprise the following components involving these entities.

- (1) *Individual reporting*: Reports are submitted by members of the affected population about specific experiences with the evaluated system.
- (2) *Aggregation for evaluation*: Reports are aggregated and interpreted over time. The goal of such aggregation is evaluation: to understand or describe the behavior of the evaluated system in a fine-grained manner.
- (3) *Evaluation-conditional action:* The aggregated evaluation supports downstream action. That is, there are evaluation outcomes where, if and when the reports are consistent with those outcomes, the mechanism administrator takes associated action.

Ensuring that reports are tied to the *evaluated system* makes it possible for the aggregation to generate specific insights about system behavior; this specificity allows for downstream action, in that it explicitly describes some (undesirable) property of the system that has emerged. The temporal



Figure 1. The components of our framework and how they interact.

component is critical—continuous evaluations make it possible to understand changing use cases and experiences over
time, and moreover, to identify problems (and take action)
quickly as they arise.

134 135

152

153

154

140 In Figure 2, we give four examples of AIR mechanisms. In 141 the first row, we show how VAERS, an existing reporting 142 system, implements the criteria of our framework. The sec-143 ond row presents a hypothetical problem setting following 144 an example in Dai et al. (2025), and the third and fourth 145 rows highlight speculative examples for applications of AIR mechanisms. As these examples indicate, AIRs can be es-147 tablished for various systems (in fact, the evaluated system 148 need not be algorithmic at all); moreover, the implementa-149 tion of the mechanism depends closely on the type of harm 150 that an AIR is designed to surface. 151

## 2.2. Current state-of-the-art in post-deployment evaluation

155 The idealized definition given in Section 2.1 excludes most 156 prominent systems that currently exist for crowdsourcing 157 and post-deployment monitoring for AI. Generally speak-158 ing, existing systems appear to fall into three categories: 159 third-party collections of real-world problems with AI sys-160 tems; general approaches to post-deployment evaluation; 161 and targeted flaw disclosure approaches like bug bounties 162 and red-teaming. We briefly clarify the relationship of aggre-163 gated individual reporting mechanisms to these approaches. 164

In the first category are several well-known incident databases and risk repositories, which include the website incidentdatabase.ai, where incident submissions are available to the public (McGregor, 2021); the OECD's AI Incidents and hazards Monitor (AIM), which is an automated system that scrapes all AI-related news headlines globally (OECD, 2023); and the MIT AI Risk Repository, which tags individual incidents from incidentdatabase.ai with additional metadata from the Risk Repository framework (Slattery et al., 2024) and is available to download as a Google sheet. While the exact implementation of each database is slightly different, the high-level commonality is that they tend to collect discrete, externally-submitted incidents. These incidents are typically news events that were related in some way to any deployed AI system, with a general focus on examples of misuse.<sup>2</sup> As a result, these collections of incidents are much broader than individual interactions that can become descriptive of model behavior at a more granular level; in fact, aggregation of these incidents happens only to the extent of tabulating the approximate frequency of these discrete news events across impact category and severity.

<sup>&</sup>lt;sup>2</sup>For instance, at time of writing, the most recent incident at incidentdatabase.ai is that the New Orleans police department used banned facial-recognition software (Incident 1075). The harm being described in this incident is about NOPD's usage of a particular system, rather than a specific instance or individual experience of that system's failure.

65 66 67	Evaluated system	Mechanism administrator	Affected population	Individual report information	<b>Evaluation condition</b> (when would downstream action be taken due to patterns in reports?)	Downstream actions
68 69 70 71 72 73	FDA-approved vaccines (deployed in real-world system as VAERS)	3 <sup>rd</sup> -party (United States CDC & FDA)	All patients who received a particular vaccine	Specific vaccine (brand and batch), specific adverse event (e.g., ) and demographic information	Elevated overall frequency of adverse event reports compared to expected baseline frequency e.g., myocarditis appears frequently for the COVID-19 vaccine.	Further investigation of specific vaccine-side effect pairs (e.g. additional research or data collection), and notification of relevant parties (e.g. published reports)
74 75 76 77 78 79	Loan allocation algorithm at Bank X (hypothetical from Dai et al. 2025)	3 <sup>rd</sup> -party (activist organization)	All loan applicants to Bank X	Demographic information and the claim of potential discrimination	Identification of a subgroup that experiences disproportionate rates of harm e.g., financially-healthy Black applicants are denied loans at a higher rate.	Gathering evidence to initiate a legal discrimination case
80 81 82 83 84	AI medical scribe product (speculative example)	2 <sup>nd</sup> -party (hospital system client)	Healthcare workers who use the tool, and their patients	Free-text notes and information about reporter; scribe text and edit history, for healthcare worker reporters	Clinically-relevant failure modes of the scribe product. e.g., AI scribe repeatedly makes errors for visits about pregnancy complications.	Feedback provided to company developing AI scribe; temporary usage guidelines given to clinicians working in (e.g.) maternal health
85 86 87 88 89	<b>ChatGPT</b> (speculative example)	1 <sup>st</sup> -party (OpenAl)	All users of the ChatGPT product	Chat transcript and free-text notes	Wide-scale safety-critical behavior. e.g., the newest model exhibits dangerously sycophantic behavior.	Rollback to prior model version; post-mortem conducted for flaws in internal evaluations.

Figure 2. Examples of how AIRs could be set up for a variety of applications. Here, we elide the corresponding aggregation methods in order to focus on the application itself; note, however, that the "evaluation conditions" are aggregate system failures rather than per-report problems.

The second category includes most current approaches to 196 post-deployment monitoring or evaluation; here, we high-197 light the approaches that are most closely related. For understanding real-world usage of Claude, Anthropic has devel-199 oped a system called Clio, which applies k-means cluster-200 ing to a fixed batch of chat transcripts, then post-hoc labels the clusters with Claude (Tamkin et al., 2024). Unlike the incident databases described above, Clio is specific to an evaluated system (Claude)-and, notably, the goal of Clio is 204 specifically to identify "unknown unknowns" in usage patterns. However, the collection of all chats is quite different 206 from individually-submitted reports that highlight problematic behavior, and thus Clio captures different insights than an AIR would. Moreover, at least according to public infor-209 mation, Clio operates on static batches of transcripts, so it is 210 unknown how or whether it can identify trends that develop 211 over time. For large language models in general, Chatbot 212 Arena has become a popular evaluation platform for ranking multiple LLMs (Chiang et al., 2024). While Chatbot 214 Arena does rely on real-time, crowdsourced feedback, the 215 goal of their evaluations is primarily to compare LLMs (by 216 generating a ranking), rather than conducting fine-grained 217 evaluations that could later inform downstream action; more-218

over, the form of feedback afforded by the platform is rather limited as a "report."

Finally, we note that *flaw disclosure* mechanisms, e.g., as outlined in Longpre et al. (2025), are an important starting point, but are insufficient for our goals. That is, an aggregated individual reporting mechanism can-and likely should-include much of the flaw disclosure machinery outlined in Longpre et al. (2025), but a flaw disclosure system by itself would not necessarily qualify as an AIR. In the context of our definition, flaw disclosure systems may sometimes satisfy the first condition (individual reporting), but do not include the other criteria (aggregation for evaluation and downstream action). Thus, our specific proposal is explicitly narrower.<sup>3</sup> This can be seen when considering concrete

190 191

193

<sup>&</sup>lt;sup>3</sup>Regarding the specific manuscript presented in Longpre et al. (2025), the contribution of our work is distinct, but complementary. We are explicitly focused on public reporting and understanding its benefits, as in Sections 3 and 4 (whereas non-expert reports are are only briefly mentioned in their work); we also emphasize methodological problems and directions for the AI research community (Section A). We refer readers to their manuscript for a thorough legal and policy overview and a taxonomy of current flaw disclosure methods, which our manuscript excludes.

220 instantiations of flaw disclosure mechanisms, such as bug 221 bounties for algorithmic problems and red-teaming. The 222 former is typically focused on resolving bugs on a per-report 223 level, rather than understanding problems and taking action 224 on an aggregate level (e.g., as discussed in (Kenway et al., 225 2022), with the bias-bounty mechanism of Globus-Harris et al. (2022) as a rare exception); the latter often focuses on 227 evaluations at the pre-deployment phase, and solicits partic-228 ipation from predefined groups rather than all members of 229 the affected population (e.g., Ahmad et al. (2025)).

230

231

232

233

234

235

236

237

238

239

We see our proposal as complementary to all of these approaches, which are useful and important parts of the AI evaluation ecosystem. AIRs, as we have defined them, enable distinct types of evaluations beyond what is already covered by these approaches—but cannot, by themselves, supersede these efforts.

# 3. Aggregated individual reporting enables actionable harm discovery

240 In this section, we argue that AIRs, as outlined in Section 241 2.1, have concrete benefits beyond what existing systems 242 can provide. Prior work, especially in human-computer 243 interaction, has documented a desire for ways to enable 244 auditing and evaluation of algorithmic systems from a wider 245 population. For example, user-driven auditing (e.g., Shen 246 et al. (2021); DeVos et al. (2022); Deng et al. (2023); Lam 247 et al. (2022)) has been studied as a means for eliciting user 248 feedback that identifies problems with an algorithmic sys-249 tem beyond centralized evaluations. Shared concerns and 250 challenges presented by these works, which have primarily 251 involved small-scale case studies, include aggregating and 252 interpreting feedback, and using feedback to drive down-253 stream action. In fact, Ojewale et al. (2025) explicitly call 254 for the development of "tools for harm discovery" and "par-255 ticipatory methods" as a pathway towards accountability. 256

Our framework seeks to formalize some these intuitions,
and each component of our definition in 2.1 plays a specific
role: individual reports enable harm discovery by surfacing
unknown unknowns, while aggregation is a prerequisite to
making them actionable.

## 3.1. Individual reports surface unknown unknowns

In domains where reporting systems are well-established, it is widely understood that reports contain useful information that may otherwise never be identified (e.g., see (Beaubien & Baker, 2002) for a decades-old review in aviation, and (Wu et al., 2002) for health). In our vision of individual reporting, however, the affected population is a much wider set of people, beyond domain experts and practitioners. What kinds of insights might we expect to see from their reports?

Prior work suggests that audit-style feedback from "end

users" draws from their prior beliefs and experiences, and can reveal clusters of shared problems (including some that were previously unknown or unaccounted for), as well as identify meaningful disagreement across users (Lam et al., 2022; DeVos et al., 2022). For concreteness, we will focus here on social media platforms like Twitter/X as as case studies for (informal) individual reporting.

The Twitter app is of course not a purpose-built reporting system, and there is therefore no aggregation mechanism dedicated to explicitly making sense of report content. However, individuals frequently take to social media to share their experiences with specific algorithms. Shen et al. (2021) term this phenomenon "everyday algorithm auditing," and show how social media platforms enable an organic, informal process for both raising awareness and validating problems raised by initial reports. In their analysis, they highlight that "everyday audits" leverage the lived experience (e.g., cultural background) and situated knowledge (e.g., application contexts) of everyday users.

Returning to the motivating example at the beginning of this paper, the analysis of Shen et al. (2021) is largely consistent with the patterns that can be seen from posts about Chat-GPT sycophancy. For example, users noted performance degradation even in quantitative tasks (Ho, 2025; Laura, 2025), reflecting situated knowledge with respect to users' expectations in specific domain applications; more serious posts highlighted scenarios where ChatGPT validated flatearth and conspiracy beliefs (fortheloveoftheworld, 2025), reflecting both lived experience and situated knowledge.

A subtly distinct pattern from that arose in reports about sycophancy, compared to trends identified by Shen et al. (2021), is reports that reflect user creativity—rather than realistic practical usage—in a way that can be thought of as informal "red teaming." For example, users showed example output responses when given prompts about math and "pickle rick" (Williawa, 2025), monologues by unsavory fictional characters (Bharath, 2025), as well as genuinely-safety critical responses elicited intentionally (Reviews, 2025; Frye, 2025).

Importantly, the sycophancy case study also reflects a limitation of relying on social media as a vehicle for "reporting." The *types* of problems that can be identified via these "reports," somewhat by definition, are those that are amenable to virality. In fact, the methodology of the taxonomy in Shen et al. (2021) also relied on prior knowledge of cases that were "high-profile" on social media. Sycophancy is a prime example of a virality-friendly problem: It appeared to affect all users across use cases and backgrounds, so that various users felt empowered to share their particular experience that reflected the underlying problem. Moreover, it was easy to produce content that illustrated sycophancy, but was also (e.g.) highly humorous or inflammatory. However, many serious model flaws are not so "clickworthy"; some failures
might only affect a small slice of users, and in ways that
cannot be constructed as a popular tweet. More systematic
approaches, therefore, are necessary.

#### 280 **3.2. Aggregation for actionability and accountability**

281

315

We now argue that, beyond individual reporting, aggrega-282 tion is necessary to make use of the information contained in 283 reports. Not all potential problems with an evaluated system 284 are inherently statistical; in many cases, however, individual 285 experiences can only be understood in the context of aggre-286 gated evaluations. For example, if the goal is to identify 287 subgroups that disproportionately experience harm, as in 288 Dai et al. (2025), then individual reports become significant 289 insofar as they can be described statistically; similarly, par-290 ticipants in DeVos et al. (2022) study mention uncertainty 291 about whether particular algorithmic behaviors ought to be 292 considered examples of harm, especially in the absence of 293 knowledge about how others experienced the system. 294

295 For areas where report databases (often called "incident 296 databases") are already established as a component of safety 297 monitoring, it is well-understood that that focusing on rectifying individual incidents is limited. Instead, the goal 299 should be instead to find high-level patterns that can inform 300 future procedural changes; in these fields, learning from 301 reports has indeed been effective for shaping future practice 302 (Macrae, 2016; Jacobsson et al., 2012; Robinson, 2019). 303

Despite this, many reporting systems do not include in-304 tentional aggregation as a core component; these systems 305 appear to have had limited impact beyond the scope of re-306 ports themselves. On the other hand, the few examples 307 of successful action from crowdsourced information have 308 involved aggregation, suggesting that aggregation is a neces-309 sary (though potentially not sufficient) condition for taking 310 high-level action from reports. In the remainder of this 311 section, we discuss several case studies of crowdsourced 312 or public reporting feedback, and the extent to which they 313 were successfully spurred concrete action. 314

316 CFPB Consumer Complaints, VAERS, and different 317 levels of actionability. The U.S. Consumer Finance Pro-318 tection Bureau (CFPB) maintains a Consumer Complaint 319 Database to which members of the public are eligible to 320 submit. The CFPB itself does not directly analyze or ad-321 dress these complaints Consumer Financial Protection Bu-322 reau (2012); instead, it acts as an intermediary, passing the 323 complaints to the relevant financial institutions and mandat-324 ing direct responses from the financial institution (Littwin, 325 2015). Though the complaints submitted by individuals do actually receive responses from the responsible parties (i.e. 327 banks), the emphasis is on directly addressing individual 328 incidents, rather than the problems that emerge when con-329

sidering them all collectively. In this way, the Consumer Complaint Database resembles the incident databases described in Section 2.2 and the flaw disclosure framework from Section 2.1.

To the best of our knowledge, insights from the CFPB complaint database as a whole has not triggered further enforcement or legislative action; this is perhaps unsurprising, given the emphasis on actionability for individual complaints, rather than at in aggregate. Several academic works have found problematic trends by analyzing complaints collectively (e.g., Bastani et al. (2019); Ayres et al. (2013); Haendler & Heimer (2021)); thus, this focus on per-complaint resolution is a design choice, not an inherent limitation of the data itself.

On the other hand, the VAERS database is continually monitored explicitly *for* aggregate-level harm. For VAERS, the event of concern is not any individual report of an adverse event. Rather, the concern is whether reports for particular vaccines occur with abnormally high frequency; examples of clinically-relevant side effects initially discovered via VAERS abound (Shimabukuro et al., 2015; Singleton et al., 1999).

Beyond initial identification, another important usage of VAERS has been post-hoc investigation of problems that were originally flagged via case study. For instance, it is now well-known that the COVID-19 vaccines induced an elevated risk of myocarditis, but only in younger men; however, in the early stages of vaccine rollout, the conversation was limited to various healthcare providers noticing, via case study, that myocarditis appeared to be a common occurrence overall (Mouch et al., 2021; Larson et al., 2021; Marshall et al., 2021). A more holistic understanding of the issue, including that myocarditis appeared to be limited to younger men, was attained only after post-hoc analysis of VAERS reports (Witberg et al., 2021; Oster et al., 2022). This, in fact, was one motivating example for the method proposed in Dai et al. (2025): if VAERS reports about myocarditis had been analyzed with a disaggregation over demographic identity, then VAERS reports could have directly confirmed the affected subgroup more quickly.

**ChatGPT sycophancy, "spiritual delusions", and the limitations of informal aggregation.** We conjecture that one reason that the ChatGPT sycophancy problem resulted in decisive change was that the Twitter timeline algorithm, and the platform's affordances of likes and retweets, served as a quasi-aggregation scheme. The brief dominance of tweets highlighting sycophancy on the timeline showed that this was a problem that impacted a wide range of users, and that there was general agreement that the behavior was problematic. In a move that is remarkably unique for tech companies, this culminated in explicit action by OpenAI to roll back the new model deployment, and to initiate some discussion
of what went wrong (OpenAI, 2025b;a). In some sense, it
was "lucky" that this particular safety problem happened
to have been friendly for virality; in general, there is no
guarantee that quasi-reports to social media are necessarily
seen, or taken seriously, by those who have control over the
evaluated system.

337 It is instructive, here, to make a comparison to some 338 approximately-contemporaneous Reddit threads, which de-339 tailed severe mental health crises due to what appear to be 340 ChatGPT-caused "spiritual delusions" (first referenced in 341 Klee (2025); the more recent Hill (2025) describes similar 342 phenomenons). Notably, many cases mentioned had been 343 ongoing for weeks, and therefore could not be entirely attributed to the late-April update. There has been no blog 345 post that explicitly addresses these problems-and, indeed, no reason to think one might be forthcoming, based on the 347 company's statements in the reported articles. 348

349 As for why this might be the case, the obvious reason is 350 that the "spiritual delusion" problem is likely more complex than the sudden increase in sycophancy, which was 351 352 easily addressed by a rollback. However, we speculate that 353 perhaps one additional factor in the lack of decisive, publi-354 cized action from OpenAI is that, overall, the "aggregation 355 mechanism" of Reddit was much less powerful than on Twit-356 ter. The conversation about long-run psychological impact 357 did surface beyond Reddit, where the posts were originally 358 made, to reported features in major national magazines. 359 However, both Reddit as a platform and reported stories as a 360 format emphasize individual-level narratives, which are less 361 compelling as indictments of systematic behavior patterns. 362

Targeted crowdsourcing and the value of statistical evidence. Finally, we discuss two systems that crowdsourced targeted surveys of specific algorithmic systems. While the scope of these surveys were narrower—meant to discover average rates of pre-specificed metrics, rather than general evaluation—they are notable because their aggregations provided concrete statistical evidence for individual experiences.

372 In 2020 and 2021, Mozilla led a study using a browser ex-373 tension called RegretsReporter, which crowdsourced infor-374 mation about the experience of Youtube recommendations 375 (McCrosky & Geurkink, 2021). The study found that that 376 recommendations from the Youtube algorithm were dispro-377 portionately responsible for serving content that users regret-378 ted seeing (and violated terms of service)-and moreover, 379 that non-English speakers were most seriously affected.<sup>4</sup> While Youtube never publicly disclosed whether specific 380 381 algorithmic changes were made in response (a weakness we

382

discuss in Section 5), the company did directly address the report in public statements (Klar, 2021; Lawler, 2021; The Next Web, 2021).

Perhaps more optimistically, *Fairfare* is a system that crowdsources information on rideshare wages, with the goal of understanding the extent to which drivers were being underpaid as well as overall average pay rates (Calacci et al., 2025). This information not only empowered drivers to organize, but also led to legislative impact that was motivated by statistics computed via Fairfare.

More broadly, statistical evidence-in contrast to, e.g., anecdotal accounts-is especially useful as a means for pressuring organizations or institutions to make change (Recht, 2025). For example, business leaders are more amenable to making decisions that appear "data-driven," rather than responsive to anecdotal experiences (H. Davenport, 2014). Statistical evidence is also treated differently in litigation contexts (see, e.g., (Espeland & Vannebo, 2007)). The question of what kinds of statistical evidence are considered acceptable will depend on the particular application context, and not all harms are inherently statistical. At the very least, however, the language of statistics can empower individuals by validating their experiences. Both RegretsReporter and Fairfare were simply *formalizing* folk intuitions that Youtube users and rideshare drivers individually already had—but the collection and aggregation of data made it impossible to dismiss those perspectives as individual experiences of one-off anomalies.

# 4. Aggregated individual reporting as a pathway towards "democratic" AI

In this section, we take a brief detour from practical benefits to discuss the normative underpinnings of our proposal.<sup>5</sup> In recent years, the recognition that modern AI systems both require and accelerate the concentration of power has spurred a flurry of research on how AI systems might be designed through quasi-democratic processes—"participatory" (Birhane et al., 2022; Delgado et al., 2023; Gilman, 2023), "pluralistic" (Dai & Fleisig, 2024; Sorensen et al., 2024a;b), "collective" (Huang et al., 2024), and so on. Methodological research in these directions focus almost exclusively on the development phase, and rarely consider how "the public" might engage with AI systems post-deployment—despite calls for increasing the power of "the public" from more conceptual work (e.g. (Feffer et al., 2023)).

The "democracy" analogy relies on a rough analogy be-

<sup>&</sup>lt;sup>4</sup>The latter is also an example of crowdsourced data finding new, non-obvious insights, as discussed in 3.1.

<sup>&</sup>lt;sup>5</sup>While of course the normative and practical cannot be fully disentangled—indeed, Estlund's core thesis about democracy is that it is practically effective exactly because it is normatively desirable, and vice versa (Estlund, 1997)—the arguments we present in this section are those that rely on conceptual foundations.

tween AI systems and governmental bodies (about which
democracy is classically theorized). For instance, these
works commonly ask, how might "democratic" inputs shape
the model spec or training objectives? Implicit in this question is an analogy to the same way that "democratic" inputs
determine the outcome of an election.

We argue that such an analogy, while not incorrect, is incomplete. A key tenet of democracy is *consent of the governed*: the ability for members of the public to express their will about governance, and a process for those expressions of will to have meaningful bearing on future outcomes.<sup>6</sup> Cru-396 cially, such consent is not a singular event; instead, it is 397 an ongoing process that, theoretically, must continue for as long as the governing body remains in power (Bertram, 399 2023; Gourevitch & Rousseau, 2018). Democratic legiti-400 macy derives not just from the fact that citizens can shape 401 the parameters of their government's future actions, but also 402 from their ability to provide input on ongoing action-to 403 revoke consent. 404

405 A critical ingredient for a "democratic" approach to AI, 406 therefore, is the ability for members of the public to col-407 lectively raise issues with systems *after* they have already 408 been deployed. In the same way that a democratic govern-409 ing body (a kratos) requires input from its citizens (demos) 410 to continue effective governance, those who operate an AI 411 system cannot fully understand the real-world behavior of 412 their system without the input of those who have interacted 413 with it. And, in the same way that the demos ought to have 414 its concerns be heard by its kratos, those who interact with 415 possibly-consequential systems also deserve for their con-416 cerns to be taken seriously, and systematically, by those who 417 build the systems.

418 In other words, "democratic AI" is not just a design problem; 419 it is an accountability problem. Our framework seeks to 420 take a step towards this by not only offering the ability for 421 individuals to provide feedback-via reporting-but also 422 providing an avenue for them to be heard. As a starting point, 423 aggregation is a lens through which (e.g.) the owners of 424 the evaluated system can understand and interpret feedback 425 ("helping the state see," in the sense of Scott (1998), by 426 surfacing details that centralized evaluations are unable to 427 understand). More theoretically, aggregated reports also 428 have the potential to develop and shape the "general will," 429 in the sense of Rousseau, from a collection of expressions 430 of "individual wills" (Gourevitch & Rousseau, 2018). 431

From a more concrete perspective, one recent work that
seeks to place AI governance mechanisms in the context
of "democratic" processes is the Democracy Levels framework of Ovadya et al. (2024). This framework rests on a
conception of "democratic processes" as involving a remit

437

(scope), constituent population, and an output decision; the hierarchy of "levels" corresponds to the extent to which the outcomes of this process have binding power over the AI system being governed.

Like the works mentioned above, the examples given in the Democracy Levels framework are primarily about predeployment rules and guidelines. However, we argue that mechanisms fo aggregated individual reporting can be seen as an almost direct stand-in their understanding of a "democratic process": the constituent population can be seen as our affected population, the remit can be thought of as our evaluated system, and the decision can be thought of the evaluation computed from the aggregated reports.<sup>7</sup> In fact, the mechanisms we propose can be direct drop-ins to their "levels"—for instance, in Level 0, aggregated reports have no impact; in Level 1, reports are accepted but do not necessarily affect the deployed system; in Level 2, the results of aggregated evaluation directly imply a default response action.

Finally, we note that, in political theory, the very notion of a "public" is a complex and contested notion. For instance, user-led audits have been proposed as pathways to developing *counterpublics*, where individuals that are marginalized with respect to (or otherwise in opposition to) the wider "public" come together to collectively build knowledge and take action (Fraser, 1990; Baik & Sridharan, 2024; Shen et al., 2021; DeVos et al., 2022). Our discussion in this section has not explicitly considered the distinction between a hegemonic "public" as opposed to marginalized "counterpublics," instead generally emphasizing the value of giving voice to a *generic* public. We leave discussion of this distinction to future work; we expect that whether AIRs do empower counterpublics will depend on the details of how practical implementation plays out.

#### 5. Discussion

This work seeks to establish *aggregated individual reporting* as a conceptual framework for post-deployment evaluation of AI systems. There are several well-founded criticisms of aggregated individual reporting; nevertheless, we see these concerns as important considerations for carefully designing and improving AIRs, rather than reasons that they should not be attempted at all.

#### 5.1. Limitations of aggregated individual reporting

We see two major categories of failure modes for AIRs: those arising from crowdsourced reporting as a data source,

<sup>&</sup>lt;sup>6</sup>This is, e.g., canonically expressed in Locke (1988).

<sup>438</sup> 439

<sup>&</sup>lt;sup>7</sup>One interesting difference is that their definition of "democratic process" assumes a one-time event that can be initiated by various discrete triggers, whereas we think of our mechanism as continuous.

440 and those arising from organizational challenges.

442

443

444

445

446

471

(1) **Reporting is a fundamentally-flawed source of data.** By design, feedback collected by AIRs is "one-sided," in the sense that, and moreover, relies on usage or adoption at scale.

447 (1a) Individual reports would be too noisy or too biased (e.g., by reports attempting to "game" the system) as to be 448 unusable. This is a serious concern, especially given the 449 450 existence of (online) crowd behaviors like review bombing (Payne, 2024), and the varying quality of complaints about 451 452 content moderation algorithms. Anecdotally, for instance, 453 user feedback about moderation is often about explicit con-454 tent.

(1b) People may not submit reports even if the mechanism technically exists—e.g., if reporting is too burdensome, or if affected populations are unaware of the option to report. Encouraging sufficient participation is also a common challenge across applications that are rely on eliciting data from the public (e.g., study recruitment and retention (Koo et al., 2005)).

463 Handling both of these concerns is partially an empirical question (in what ways would a reporting system for evalu-464 465 ating AI induce different behaviors? What kinds of report 466 affordances affect incentives?) that can only be understood 467 by analyzing a real-world implementation. On the other hand, as mentioned in A.2, there are also various research 468 469 communities that can and should contribute towards these 470 problems.

(2) Organizational challenges may complicate the pathway to downstream action or accountability. One key premise of our proposal is that reports can be empowering because they can effect action, rather than languishing in an online database. However, for this to happen, there are organizational and institutional problems that must be addressed beyond technical and methodological challenges.

(2a) Model developers may not be incentive-aligned as 1st-480 party operators. Even if problems with the evaluated system 481 are identified, the system developer is not necessarily bound 482 to address them. Though Mozilla's Youtube Regrets study 483 was discussed earlier as a positive example of aggregation 484 (McCrosky & Geurkink, 2021), Youtube never explicitly ac-485 knowledged the study as influencing specific choices about 486 their recommendation algorithm. On the other hand, a first-487 party administrator could evade accountability by avoiding 488 disclosure of findings from the AIR system. 489

(2b) Running an AIR mechanism as a 3rd-party may be
unsustainable. Many third-party audits are foundationfunded (e.g., RegretsReporter by Mozilla), and the path
to long-term financial viability is uncertain, which means

that many 3rd-party auditors simply cease to exist. For example, the previously-lauded UberCheats browser extension (Marshall, 2021) no longer exists; the system that became Fairfare (Calacci et al., 2025) was originally known as the DriversSeat app (Driver's Seat Cooperative, 2024), but it is unclear whether prior data from DriversSeat was ever used.

These are critical concerns about ideal patterns of implementation, especially when considering which institutions ought to play what roles. We cannot guarantee *a priori* that these failure modes can be avoided, but we hope that by emphasizing the role of organizational factors in A.1, we can encourage intentional decisionmaking in this regard.

#### 5.2. Calls to action

Despite these limitations, we believe that aggregated individual reporting, as we have defined in this work, is a natural component of the post-deployment evaluation ecosystem. Recent events—and ongoing research—have illustrated that individuals have unique contributions to understanding the contours of AI system behavior. The constellations of individual experience are an invaluable resource not just for model development, but for understanding potentiallyunintended societal consequences of already-deployed systems.

Our main call to action, therefore, is for AIRs to be built and the data collected by them to be analyzed. *Academic researchers* should develop the methodological innovations and empirical analyses necessary for effective implementation; this is an interdisciplinary challenge that spans several subfields of computer science (and beyond), including not just AI/ML but also statistics, human-computer interaction, and law and policy. *Activists and industry practitioners* should begin exploring development of these systems. *Policymakers* should work in collaboration with the aforementioned stakeholders to understand what kinds of legal or policy leverage may be useful.

To translate the AIR approach from a hypothetical strategy to reality, it is essential for the area to mature; addressing the wide range of design and methodological issues outlined in this position is just the first step in coalescing a community to make progress towards this goal. But, at the same time, it is impossible to wait for all the hypothetical kinks to be smoothed: the success of AIRs is a fundamentally empirical question. We believe that it is worth trying to find out.

#### References

495

496

497

498

499

500

505

506

507

508

509

511

512

513

514

520

525

535

536

537

538

- Agostini, G., Pierson, E., and Garg, N. A Bayesian Spatial Model to Correct Under-Reporting in Urban Crowdsourcing. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 21888–21896, 2024.
- 501 Ahmad, L., Agarwal, S., Lampe, M., and Mishkin, P. Ope-502 nai's approach to external red teaming for ai models and 503 systems. arXiv preprint arXiv:2503.16431, 2025. 504
- Assembly, U. G. Seizing the opportunities of safe, secure and trustworthy artificial intelligence systems for sustainable development, 2024. URL https://documents.un.org/doc/undoc/ ltd/n24/065/92/pdf/n2406592.pdf. 510
  - Ayres, I., Lingwall, J., and Steinway, S. Skeletons in the database: An early analysis of the cfpb's consumer complaints. Fordham J. Corp. & Fin. L., 19:343, 2013.
- 515 Baik, J. and Sridharan, H. Civil rights audits as counter-516 public strategy: articulating the responsibility and failure 517 to care for marginalized communities in platform gover-518 nance. Information, Communication & Society, 27(5): 519 836-855, 2024.
- Bastani, K., Namavari, H., and Shaffer, J. Latent dirichlet 521 522 allocation (lda) for topic modeling of the cfpb consumer 523 complaints. Expert Systems with Applications, 127:256-524 271, 2019.
- Beaubien, J. M. and Baker, D. P. A review of selected 526 aviation human factors taxonomies, accident/incident re-527 porting systems and data collection tools. International 528 Journal of Applied Aviation Studies, 2(2):11–36, 2002. 529
- 530 Bertram, C. Jean Jacques Rousseau. In Zalta, E. N. and 531 Nodelman, U. (eds.), The Stanford Encyclopedia of Phi-532 losophy. Metaphysics Research Lab, Stanford University, 533 Summer 2023 edition, 2023. 534

Bharath, S. Post on Χ. https: //x.com/Siddharth87/status/ 1916999455146185022, April 2025.

- 539 Biden, J. R. Executive order on the safe, secure, and trust-540 worthy development and use of artificial intelligence. 541 2023. 542
- 543 Birhane, A., Isaac, W., Prabhakaran, V., Diaz, M., Elish, 544 M. C., Gabriel, I., and Mohamed, S. Power to the people? 545 Opportunities and challenges for participatory AI. In 546 Proceedings of the 2nd ACM Conference on Equity and 547 Access in Algorithms, Mechanisms, and Optimization, pp. 548 1-8, 2022. 549

- Calacci, D., Rao, V. N., Dalal, S., Di, C., Pua, K.-W., Schwartz, A., Spitzberg, D., and Monroy-Hernández, A. Fairfare: A tool for crowdsourcing rideshare data to empower labor organizers. arXiv preprint arXiv:2502.11273, 2025.
- Chausson, S., Fourcade, M., Harding, D. J., Ross, B., and Renard, G. The insight-inference loop: Efficient text classification via natural language inference and threshold-tuning. Sociological Methods & Research, pp. 00491241251326819, 2025.
- Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhu, B., Zhang, H., Jordan, M., Gonzalez, J. E., et al. Chatbot arena: An open platform for evaluating llms by human preference. In Forty-first International Conference on Machine Learning, 2024.
- Consumer Financial Protection Bureau. Consumer complaint database, June 2012. URL https://www.consumerfinance.gov/ data-research/consumer-complaints/.
- Dai, J. and Fleisig, E. Mapping social choice theory to RLHF. arXiv preprint arXiv:2404.13038, 2024.
- Dai, J., Gradu, P., Raji, I. D., and Recht, B. From individual experience to collective evidence: A reporting-based framework for identifying systemic harms. International Conference on Machine Learning, 2025.
- Delgado, F., Yang, S., Madaio, M., and Yang, Q. The participatory turn in AI design: Theoretical foundations and the current state of practice. In Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization, pp. 1–23, 2023.
- Deng, W. H., Guo, B., Devrio, A., Shen, H., Eslami, M., and Holstein, K. Understanding practices, challenges, and opportunities for user-engaged algorithm auditing in industry practice. In Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, pp. 1-18, 2023.
- DeVos, A., Dhabalia, A., Shen, H., Holstein, K., and Eslami, M. Toward user-driven algorithm auditing: Investigating users' strategies for uncovering harmful algorithmic behavior. In Proceedings of the 2022 CHI conference on human factors in computing systems, pp. 1–19, 2022.
- Driver's Seat Cooperative. Driver's seat cooperative. https://www.driversseat.co/,2024.
- Espeland, W. N. and Vannebo, B. I. Accountability, quantification, and law. Annu. Rev. Law Soc. Sci., 3(1):21-43, 2007.

- Estlund, D. Beyond fairness and deliberation: The epistemic dimension of democratic authority. *Deliberative democracy: Essays on reason and politics*, 173:204, 1997.
- European Parliament. Regulation (eu) 2024/1689 of the
  european parliament and of the council of 13 june 2024
  laying down harmonised rules on artificial intelligence
  (artificial intelligence act), 2024. Chapter IX, Articles 85
  and 87.
- Feffer, M., Skirpan, M., Lipton, Z., and Heidari, H. From preference elicitation to participatory ml: A critical survey & guidelines for future research. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 38–48, 2023.
- 565 fortheloveoftheworld. This new update is unacceptable and absolutely terrifying. https://www. reddit.com/r/OpenAI/comments/lkasjmr/ this\_new\_update\_is\_unacceptable\_and\_ absolutely/, April 2025.
- 571 Fraser, N. Rethinking the public sphere: A contribution to
  572 the critique of actually existing democracy'. *Social Text*,
  573 (25/26):56–80, 1990.

570

574

575

- Frye. Post on x. https://x.com/\_\_\_frye/ status/1916346474893656572, April 2025.
- 577 Ghosh, A. and McAfee, P. Incentivizing high-quality user578 generated content. In *Proceedings of the 20th interna-*579 *tional conference on World wide web*, pp. 137–146, 2011.
- 581 Gilman, M. E. Democratizing AI: Principles for meaningful
  public participation. *Data & Society*, 2023.
- Globus-Harris, I., Kearns, M., and Roth, A. An algorithmic framework for bias bounties. In *Proceedings of the* 2022 ACM Conference on Fairness, Accountability, and Transparency, pp. 1106–1124, 2022.
- Globus-Harris, I., Harrison, D., Kearns, M., Perona, P., and Roth, A. Diversified ensembling: An experiment in crowdsourced machine learning. In *Proceedings of the* 2024 ACM Conference on Fairness, Accountability, and Transparency, pp. 529–545, 2024.
- Gourevitch, V. and Rousseau, J.-J. *Rousseau: the Social Contract and other later political writings*. Cambridge
  University Press, 2018.
- H. Davenport, T. How strategists use "big data" to support internal business decisions, discovery and production. *Strategy & leadership*, 42(4):45–50, 2014.
- Haendler, C. and Heimer, R. The financial restitution gap in consumer finance: Insights from complaints filed with the cfpb. *Available at SSRN*, 3766485, 2021.

- Hill, K. They asked chatgpt questions. the answers them spiraling. The New sent York Times, 2025. URL https://www. nytimes.com/2025/06/13/technology/ chatgpt-ai-chatbots-conspiracies. html.
- Ho, E. Post on x. https://x.com/Eddie\_Ho2025/ status/1913596114156028013, April 2025.
- Hu, X., Jagadeesan, M., Jordan, M. I., and Steinhardt, J. Incentivizing high-quality content in online recommender systems. *arXiv preprint arXiv:2306.07479*, 2023.
- Huang, S., Siddarth, D., Lovitt, L., Liao, T. I., Durmus, E., Tamkin, A., and Ganguli, D. Collective Constitutional AI: Aligning a Language Model with Public Input. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1395–1417, 2024.
- Jacobsson, A., Ek, Å., and Akselsson, R. Learning from incidents–a method for assessing the effectiveness of the learning cycle. *Journal of Loss Prevention in the Process Industries*, 25(3):561–570, 2012.
- Kalluri, P. Don't ask if artificial intelligence is good or fair, ask how it shifts power. *Nature*, 583(7815):169–169, 2020.
- Kenway, J., François, C., Costanza-Chock, S., Raji, I. D., and Buolamwini, J. Bug bounties for algorithmic harms? lessons from cybersecurity vulnerability disclosure for algorithmic harms discovery, disclosure, and redress. *Algorithmic Justice League*, 2022.
- Klar, R. Youtube algorithm keeps recommending 'regrettable' videos, study finds. The Hill, July 2021. URL https:// thehill.com/policy/technology/ 561788-youtube-algorithm-keeps-recommending-regre
- Klee, M. Ai spiritual delusions are destroying human relationships. *Rolling Stone*, May 2025. URL https://www.rollingstone. com/culture/culture-features/ ai-spiritual-delusions-destroying-human-relations
- Koo, M., Skinner, H., et al. Challenges of internet recruitment: a case study with disappointing results. *Journal of medical Internet research*, 7(1):e126, 2005.
- Lam, M. S., Gordon, M. L., Metaxa, D., Hancock, J. T., Landay, J. A., and Bernstein, M. S. End-user audits: A system empowering communities to lead large-scale investigations of harmful algorithmic behavior. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW2):1–34, 2022.

605 606 607 608	Larson, K. F., Ammirati, E., Adler, E. D., Cooper Jr, L. T., Hong, K. N., Saponara, G., Couri, D., Cereda, A., Proco- pio, A., Cavalotti, C., et al. Myocarditis after bnt162b2 and mrna-1273 vaccination. <i>Circulation</i> , 144(6):506–508, 2021	algorithm. Technical report, Mozilla Foundation, July 2021. URL https://assets.mofoprod.net/ network/documents/Mozilla_YouTube_ Regrets_Report.pdf.
609 610 611 612 613	<pre>2021. Laura. Post on x. https://x.com/IndieLauraSDG/ status/1913102627837206587, April 2025.</pre>	McGregor, S. Preventing repeated real world ai failures by cataloging incidents: The ai incident database. In <i>Proceedings of the AAAI Conference on Artificial Intelli-</i> <i>gence</i> , volume 35, pp. 15458–15463, 2021.
614 615 616 617 618	Lawler, R. Mozilla's regretsreporter data shows youtube keeps recommending harmful videos. <i>The Verge</i> , July 2021. URL https: //www.theverge.com/2021/7/7/22567640/ youtube-algorithm-suggestions-radicali	<ul> <li>Mouch, S. A., Roguin, A., Hellou, E., Ishai, A., Shoshan, U., Mahamid, L., Zoabi, M., Aisman, M., Goldschmid, N., and Yanay, N. B. Myocarditis following COVID-19</li> <li>ZatimRNA vaccination. vaccine, 39(29):3790–3793, 2021.</li> </ul>
619 620 621	<ul><li>Littwin, A. Why process complaints: Then and now consumer. <i>Temple Law Review</i>, 87:895, 2015.</li><li>Liu, Z. and Garg, N. Equity in resident crowdsourcing:</li></ul>	<ul> <li>Movva, R., Peng, K., Garg, N., Kleinberg, J., and Pierson,</li> <li>E. Sparse autoencoders for hypothesis generation. <i>arXiv</i> preprint arXiv:2502.04382, 2025.</li> </ul>
622 623 624 625	Measuring under-reporting without ground truth data. In Proceedings of the 23rd ACM Conference on Economics and Computation, pp. 1016–1017, 2022.	OECD. Oecd ai incidents and hazards monitor (aim), 2023. URL https://oecd.ai/en/incidents.
626 627 628 629	Liu, Z., Bhandaram, U., and Garg, N. Quantifying spa- tial under-reporting disparities in resident crowdsourcing. <i>Nature Computational Science</i> , 4(1):57–65, 2024.	Ojewale, V., Steed, R., Vecchione, B., Birhane, A., and Raji, I. D. Towards ai accountability infrastructure: Gaps and opportunities in ai audit tooling. In <i>Proceedings of the</i> 2025 CHI Conference on Human Factors in Computing
630 631	Locke, J. <i>Two Treatises of Government</i> . Cambridge University Press, 1988. Originally published in 1689.	<i>Systems</i> , pp. 1–29, 2025.
632 633 634 635	Longpre, S., Klyman, K., Appel, R. E., Kapoor, S., Bom- masani, R., Sahar, M., McGregor, S., Ghosh, A., Blili- Hamelin, B., Butters, N., et al. In-house evaluation is	OpenAI. Expanding on what we missed with sycophancy. OpenAI Blog, May 2025a. URL https://openai. com/index/expanding-on-sycophancy/.
636 637 638	not enough: Towards robust third-party flaw disclosure for general-purpose ai. <i>arXiv preprint arXiv:2503.16861</i> , 2025.	OpenAI. Sycophancy in gpt-40: what hap- pened and what we're doing about it, April 2025b. URL https://openai.com/index/
639 640	Macrae, C. The problem with incident reporting. <i>BMJ</i> auality & safety 25(2):71–75 2016	Octor M.E. Show D.K. Su. L.D. Coo J. Crooch C. D.
641 642 643 644 645	<ul> <li>Malinen, S. Understanding user participation in online communities: A systematic literature review of empirical studies. <i>Computers in human behavior</i>, 46:228–238, 2015.</li> </ul>	Oster, M. E., Shay, D. K., Su, J. K., Gee, J., Creech, C. B., Broder, K. R., Edwards, K., Soslow, J. H., Dendy, J. M., Schlaudecker, E., et al. Myocarditis cases reported after mrna-based covid-19 vaccination in the us from december 2020 to august 2021. <i>Jama</i> , 327(4):331–340, 2022.
646 647 648 649 650	Marshall, A. Gig workers gather their own data to check the algorithm's math. <i>WIRED</i> , February 2021. URL https://www.wired.com/story/ gig-workers-gather-data-check-algorith	Ovadya, A., Thorburn, L., Redman, K., Devine, F., Milli, S., Revel, M., Konya, A., and Kasirzadeh, A. Toward democracy levels for ai. <i>arXiv preprint arXiv:2411.09222</i> , m-mattr/:
651 652 653 654	Marshall, M., Ferguson, I. D., Lewis, P., Jaggi, P., Gagliardo, C., Collins, J. S., Shaughnessy, R., Caron, R., Fuss, C., Corbin, K. J. E., et al. Symptomatic acute myocarditis in	<ul><li>Payne, W. B. Review bombing the platformed city: Contested political speech in online local reviews. <i>Big Data &amp; Society</i>, 11(3):20539517241275879, 2024.</li></ul>
655 656	<i>Pediatrics</i> , 148(3), 2021.	Raji, I. D., Xu, P., Honigsberg, C., and Ho, D. Outsider
657 658 659	McCrosky, J. and Geurkink, B. Youtube regrets: A crowd- sourced investigation into youtube's recommendation	governance. In <i>Proceedings of the 2022 AAAI/ACM Con-</i> ference on AI, Ethics, and Society, pp. 557–571, 2022.
	1	12

Rao, V. N., Agarwal, E., Dalal, S., Calacci, D., and Monroy-661 Hernández, A. Quallm: An llm-based framework to 662 extract quantitative insights from online forums. arXiv 663 preprint arXiv:2405.05345, 2024.

664

665

666

667

668

669

670

671

672

673

674

677

691

697

Recht, B. A bureaucratic theory of statistics. arXiv preprint arXiv:2501.03457, 2025.

Reviews, P. https:// Post on x. x.com/ReviewsPossum/status/ 1917292836082397355, April 2025.

- Robinson, S. Temporal topic modeling applied to aviation safety reports: A subject matter expert review. Safety science, 116:275-286, 2019.
- 675 Scott, J. Seeing Like a State: How Certain Schemes to 676 Improve the Human Condition Have Failed. The Institution for Social and Policy Studies. Yale University Press, 678 1998. ISBN 9780300078152. URL https://books. 679 google.com/books?id=PqcPCgsr2u0C. 680
- 681 Shen, H., DeVos, A., Eslami, M., and Holstein, K. Everyday 682 algorithm auditing: Understanding the power of everyday 683 users in surfacing harmful algorithmic behaviors. Pro-684 ceedings of the ACM on Human-Computer Interaction, 5 685 (CSCW2):1-29, 2021. 686

687 Shimabukuro, T. T., Nguyen, M., Martin, D., and DeSte-688 fano, F. Safety monitoring in the vaccine adverse event 689 reporting system (VAERS). Vaccine, 33(36):4398-4405, 690 2015.

Singleton, J. A., Lloyd, J. C., Mootrey, G. T., Salive, M. E., 692 Chen, R. T., Ellenberg, S., Rastogi, S., Krueger, C., Braun, 693 M., Wise, R., et al. An overview of the vaccine adverse 694 event reporting system (vaers) as a surveillance system. 695 Vaccine, 17(22):2908-2917, 1999. 696

Slattery, P., Saeri, A. K., Grundy, E. A., Graham, J., Noetel, 698 M., Uuk, R., Dao, J., Pour, S., Casper, S., and Thompson, 699 N. The ai risk repository: A comprehensive meta-review, 700 database, and taxonomy of risks from artificial intelligence. arXiv preprint arXiv:2408.12622, 2024.

- Sorensen, T., Jiang, L., Hwang, J. D., Levine, S., Pyatkin, 704 V., West, P., Dziri, N., Lu, X., Rao, K., Bhagavatula, C., 705 et al. Value kaleidoscope: Engaging AI with pluralistic 706 human values, rights, and duties. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, 708 pp. 19937-19947, 2024a. 709
- 710 Sorensen, T., Moore, J., Fisher, J., Gordon, M., Mireshghal-711 lah, N., Rytting, C. M., Ye, A., Jiang, L., Lu, X., Dziri, N., 712 et al. A roadmap to pluralistic alignment. arXiv preprint 713 arXiv:2402.05070, 2024b. 714

- Tamkin, A., McCain, M., Handa, K., Durmus, E., Lovitt, L., Rathi, A., Huang, S., Mountfield, A., Hong, J., Ritchie, S., et al. Clio: Privacy-preserving insights into real-world ai use. arXiv preprint arXiv:2412.13678, 2024.
- The Next Web. Youtube recommends videos that violate the platform's own policies, study finds. The Next Web, July 2021. URL https://thenextweb.com/news/ youtube-algorithm-recommends-videos-that-violate-
- Vovk, V. and Wang, R. E-values: Calibration, combination and applications. The Annals of Statistics, 49(3):1736-1754, 2021.
- Wald, A. and Wolfowitz, J. Optimum character of the sequential probability ratio test. The Annals of Mathematical Statistics, pp. 326-339, 1948.
- Williawa. Post on x. https://x.com/williawa/ status/1916935712550551743, April 2025.
- Witberg, G., Barda, N., Hoss, S., Richter, I., Wiessman, M., Aviv, Y., Grinberg, T., Auster, O., Dagan, N., Balicer, R. D., et al. Myocarditis after Covid-19 vaccination in a large health care organization. New England Journal of Medicine, 385(23):2132-2139, 2021.
- Wu, A. W., Pronovost, P., and Morlock, L. Icu incident reporting systems. Journal of critical care, 17(2):86-94, 2002.

#### A. Towards implementation: practical details and open research questions

We now turn to discussing design decisions that must be made, either implicitly or explicitly, in the implementation of any AIR; each step also naturally gives rise to various multidisciplinary research questions. In Figures 3 and 4, we outline the relevant questions in detail.

#### A.1. Concrete design decisions

Here, we overview the high-level categories of design decisions and some implications of those choices. While these decisions are interdependent—reporting affordances also affect what aggregation methods would be effective, as well as what types of evaluations are available—there are a wide range of ways in which these decisions could be made, depending on context.

Within this section, we use Dai et al. (2025) as a running example of a methodological proposal that is largely consistent with our framework; this work proposes individual reporting for post-deployment fairness auditing, and provides algorithms that are specific to this task.

**Organizational decisions.** The first core set of decisions that must be made are organizational: what entity will take the role of *mechanism administrator*, and what relationship will it have to the *evaluated system*? This decision affects what kinds of problems the mechanism hopes to identify (i.e., the end-goal of the evaluation) as well as the nature of available downstream action.

A first-party administrator will have the fullest information about the evaluated system (e.g., when particular feature rollouts or updates were made), and be able to quickly gather additional data that may become relevant in order to contextualize information raised by reports. Since the first-party administrator has control over the evaluated system, the organization can also directly make changes in reponse to evaluation results. However, a first-party administrator may also inadvertently restrict the scope of the system, or intentionally choose to ignore the evaluation. On the other end of the spectrum, a third-party administrator would have nearly no additional information beyond what is included among reports, and cannot directly update or improve the system. However, such an organization could bring external leverage to the evaluated system, e.g. via media or legal pressure.

Dai et al. (2025) is developed under an assumption that the goal would be to identify the subgroups that are harmed, but, as a primarily methodological work, it does not specify what organizational entity would be the administrator, or what downstream action might look like.

**Reporting affordances.** What information is included in a report, and how do reporters experience the process of reporting? Examples of potential report formats can be seen in the fourth column of Figure 2, as well as the the flaw disclosure worksheet in Longpre et al. (2025). For the algorithm proposed by Dai et al. (2025), the only information collected in a report is demographic information about the reporter; however, one could naturally imagine that reports could include more information depending on the application, such as medical history (for a vaccine or pharmaceutical system) or financial background information (for a loan allocation system).

**Reporting behavior.** Due to the nature of reporting data it is, intrinsically, essential to understand reporting behavior: what factors affect the decision to submit a report, and, crucially, in what ways might reports be correlated to the target of evaluation? Do different subpopulations report at different rates? Do different types of issues lead to different reporting behaviors? In Dai et al. (2025), the choice is to commit to a set of (quantitative) assumptions about the extent to which reporting rates can vary, and incorporate those assumptions into the design of the algorithm.

Aggregation method and evaluation condition. Given the affordances offered to reporters, what method will be ued to interpret them over time—that is, how are evaluation results computed? Given that method, what specific results would define the evaluation condition (i.e. trigger for downstream action)? Methods should be sequential, or at least explicitly consider a temporal component (e.g., via sequences of batched data). This is because AIRs should be accessible to reporters at any time (rather than within the scope of a centralized study with a defined start and end date), as well as the classic distribution shift problem: users' needs in relation to an evaluated system will change over time, as will the system itself. In Dai et al. (2025), the method is formalized as a sequential hypothesis test with type-I error control. Therefore, the evaluation result that triggers action is exactly when the "null hypothesis" of no overrepresented subgroups can be rejected at level  $\alpha$ .

<ul> <li>First-party: same org (e.g. OpenAl for ChatGPT)</li> <li>Second-party: user of evaluated system (e.g., hospital system for Al</li> </ul>	What arrangement is "optimal"
<ul> <li>First-party: same org (e.g. OpenAl for ChatGPT)</li> <li>Second-party: user of evaluated system (e.g., hospital system for Al</li> </ul>	• What arrangement is "optimal"
<ul><li>scribe product)</li><li>Third-party: External org (e.g., government or activist nonprofit)</li></ul>	<ul> <li>Incentive-compatible?</li> <li>How do individuals within thesorganizations conceive of pathto to impact?</li> </ul>
<ul> <li>System users (e.g., ChatGPT users)</li> <li>System users and those close to them (e.g. friends/family of ChatGPT users)</li> <li>System users and non-users affected by system usage (e.g., healthcare workers and patients)</li> </ul>	How does the inclusion of diffe "user roles" affect the substance report content?
(Reports are submitted by members of the affected population about specifi	c experiences with the evaluated sys
<ul> <li>Reporting option directly available at or after each system interaction (e.g., "submit report" available in UI, or sent as part of follow-up)</li> <li>Advertisements on social media about the opportunity to submit reports</li> </ul>	<ul> <li>What pathway is the most effect</li> <li>How do publicity methods affer who submits reports and why?</li> </ul>
<ul> <li>Submitted reports are taken "as is" with the understanding that reporting</li> <li>Assumptions are made about reporting behavior (e.g., that heterogeneity in reporting rates is not too extreme)</li> <li>Side information (outside of submitted reports) is sought in order to characterize reporting behavior (e.g., choosing some subset of reports to "verify")</li> </ul>	<ul> <li>Who is more likely to report, ar why? What topics/types of prob are more likely to generate report and why?</li> <li>How can we model and detect disparate rates of reporting for various reasons? How can we inferences that account for or a robust to uncertainties in report</li> </ul>
<ul> <li>About specific one-off interactions (e.g., "I said X and the model responded Y")</li> <li>About longer-run series of interactions (e.g., "over the last 2 weeks, the model has been telling my partner that they are a 'chosen one'")</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain additional contextual information about the reporter or impacted party (e.g., demographics, context on usage, location, timestamp, etc.)</li> <li>May or may not include information about believed severity, justification, or proposed solution (e.g., "this was a minor problem that could worsen"; "this was problematic due to X"; "it would have been better if Y happened instead")</li> </ul>	<ul> <li>How do different affordances in report format enable reports about different types of harm?</li> <li>How can we quantitatively prodifferent types of information, including rich unstructured data like free text, and incorporate the into a quantitative (possibly-statistical) aggregation?</li> </ul>
<ul> <li>No additional incentive for reporters</li> <li>Financial compensation for reports that turn out to reflect what is later deemed to be a 'true' problem</li> <li>Report affordances are designed to shape reporting behavior (e.g., intermediate report summaries are made visible to the public, so that reporters may be interested in contributing to the conversation about currently-leading concerns)</li> </ul>	<ul> <li>What motivates potential report to actually submit reports?</li> <li>Are there concerns about misaligned incentives that migh affect the 'quality' of reports?</li> <li>Are there feedback loops that m arise?</li> </ul>
Figure 3. Organizational and interaction-focused question	ons.
	<ul> <li>System users (e.g., ChatGPT users)</li> <li>System users and those close to them (e.g. friends/family of ChatGPT users)</li> <li>System users and non-users affected by system usage (e.g., healthcare workers and patients)</li> <li><i>(Reports are submitted by members of the affected population about specifi</i> (e.g., "submit report" available at or after each system interaction (e.g., "submit report" available in UI, or sent as part of follow-up)</li> <li>Advertisements on social media about the opportunity to submit reports</li> <li>Submitted reports are taken "as is" with the understanding that reporting</li> <li>Assumptions are made about reporting behavior (e.g., that heterogeneity in reporting rates is not too extreme)</li> <li>Side information (outside of submitted reports) is sought in order to characterize reporting behavior (e.g., choosing some subset of reports to "verify")</li> <li>About specific one-off interactions (e.g., "I said X and the model responded Y")</li> <li>About specific one-off interactions (e.g., "over the last 2 weeks, the model has been telling my partner that they are a 'chosen one'")</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "state" to completely reproduce the problem (e.g., real chat transcript may or may not be available)</li> <li>May or may not contain sufficient "st</li></ul>

A failed hypothesis test might, for instance, prompt further inquiry into the reported harm for a particular subgroup.

	Options and examples	Example research questions				
Online aggregatior	(Reports are aggregated and interpreted over time; the goal is describing sys	stem behavior in a fine-grained manner)				
Evaluation trigger (what evaluation seeks to identify, and what would trigger downstream action)	<ul> <li>Set up with specific type of harm in mind; try to identify the specific harm as soon as possible if it does arise</li> <li>Set up with a more unstructured approach; try to identify important trends that emerge over time, and allow</li> </ul>	<ul> <li>What methods can successfully achieve these goals in an online of quasi-online manner? How do the adapt to or handle changes over time?</li> <li>What methods can adapt to riche "hypotheses"?</li> <li>How can we balance batching da with handling true sequentiality?</li> </ul>				
Type of evidence as an output of aggregation	<ul> <li>A rigorous statistical guarantee is required (e.g., "with probability 1-α, the findings of the aggregation are significant")</li> <li>More ad-hoc interpretation is sufficient (e.g., looking at the output of a clustering algorithm)</li> </ul>	<ul> <li>What are the relative strengths an weaknesses of various types of evidence for the purpose of motivating downstream action?</li> <li>What are the tradeoffs involved ir pursuing more vs. less "rigorous" outcomes? What methods strike a balance between data efficiency validity of conclusion?</li> </ul>				
Mechanism for action (Some evaluation results may trigger action; the mechanism administrator is responsible for doing so.)						
Types of action that can be taken once evaluation trigger is reached	<ul> <li>First-party administrators: Rollback to prior model version (e.g., if new problems suddenly arise after new deployment)</li> <li>Second-party administrators: Revisit usage policy (e.g., internal guidelines for when a tool should be used or not) and/or provide feedback to owner of evaluated system</li> <li>Third-party administrators: Build public pressure (e.g., via media) and/or action towards accountability (e.g., legal case against first-party)</li> <li>All administrators: may use the evaluation trigger as a starting point for further investigation (e.g., additional data collection)</li> </ul>	<ul> <li>What kinds of evidence can compaction from an external third party</li> <li>How can closed systems be pressured to admit a reporting system?</li> </ul>				
	Public notice of problem identified & changes made	To what extent does, and should, the reporting mechanism encoura				

Finally, we highlight some examples of existing lines of work within the AI research community—beyond auditing and evaluation—that can effectively inform these design questions.<sup>8</sup> While these examples are non-exhaustive, we hope to illustrate that AIRs can benefit from a wide range of methodological and disciplinary perspectives.

**Interaction design.** From prior work, we know that the affordances that are available to users to share feedback affect the content they share (e.g., social media platforms and virality-friendly content; see also discussion in (DeVos et al., 2022) about the importance of platform affordances). What kinds of designs can encourage the most effective reporting feedback (e.g., can facilitate long term engagement, as discussed in (Deng et al., 2023))? Are checklists, as formalized in Longpre et al. (2025), sufficient? Should intermediate evaluation results be made public, and if so, how? E.g., if individuals see that others have reported similar problems, does that encourage them to submit their own reports?

 <sup>&</sup>lt;sup>877</sup>
 <sup>8</sup>While there are, of course, domain experts from other fields that have insights on the design and implementation of reporting systems, we focus here on computer science-adjacent subfields.

**Reporting rates, incentives, and behavior.** Empirical studies of existing reporting systems have shown that different subgroups often report different types of problems at different rates (e.g. (Agostini et al., 2024; Liu & Garg, 2022; Liu et al., 2024)); at the same time, these works also suggest the potential for statistical methods that can estimate or correct for varying reporting rates. As the crowdsourcing experiments of Globus-Harris et al. (2024) suggest, participants can occasionally be adversarial, and the details of mechanism implementation should be careful to incentivize "good" reporting behavior. Prior work in motivations and incentives of users on online platforms also suggest avenues for study, both qualitatively (e.g., Malinen (2015)) and via theoretical models (e.g., Ghosh & McAfee (2011); Hu et al. (2023)).

Making sense of unstructured text. Prior work has highlighted the challenge of bridging qualitative and quantitative
insights (e.g., DeVos et al. (2022); Deng et al. (2023)). Recent methodological work to this end leverages developments in
LLMs (e.g., Rao et al. (2024); Movva et al. (2025); Tamkin et al. (2024); Chausson et al. (2025)), which may prove fruitful.
However, we note that even more classical NLP approaches (e.g. as in Robinson (2019); Ayres et al. (2013); Bastani et al.
(2019)) have shown to be effective when processing existing (natural-language) report data.

Sequential analysis and online decisionmaking. Statistics has studied sequential testing for decades (e.g., Wald's SPRT (Wald & Wolfowitz, 1948)). More recently, e-values (e.g., (Vovk & Wang, 2021)) have emerged as a popular framework for sequential approaches; however, as (Dai et al., 2025) note, nontrivial extensions for settings more complex than a single hypothesis are open technical problems. Thus, it may be fruitful to explore methods that loosen the stringency of a true sequential hypothesis test, and/or that seek to reconcile the statistical approach with the text-based methods discussed above. In particular, insights from the rich literature in multi-armed bandits and other approaches to online optimization and decisionmaking may offer methodological contributions to handling sequentiality.