# Going public: the role of public participation approaches in commercial AI labs

Lara Groves, Aidan Peppin, Andrew Strait and Jenny Brennan

Ada Lovelace Institute, London, UK

## Abstract

In recent years, discussions of responsible AI practices have seen growing support for 'participatory AI' approaches, intended to involve members of the public in the design and development of AI systems. Prior research has identified a lack of standardised methods or approaches for how to use participatory approaches in the AI development process. At present, there is a dearth of evidence on attitudes to and approaches for participation in the sites driving major AI developments: commercial AI labs. Through 12 semi-structured interviews with industry practitioners and subject-matter experts, this paper explores how commercial AI labs understand participatory AI approaches and the obstacles they have faced implementing these practices in the development of AI systems and research. We find that while interviewees view participation as a normative project that helps achieve 'societally beneficial' AI systems, practitioners face numerous barriers to embedding participatory approaches in their companies: participation is expensive and resource intensive, it is 'atomised' within companies, there is concern about exploitation, there is no incentive to be transparent about its adoption, and it is complicated by a lack of clear context. These barriers result in a piecemeal approach to participation that confers no decision-making power to participants and has little ongoing impact for AI labs. This paper's contribution is to provide novel empirical research on the implementation of public participation in commercial AI labs, and shed light on the current challenges of using participatory approaches in this context.

## 1 Introduction

Artificial intelligence research and technology continues to proliferate widely, presenting substantial opportunities but also considerable ethical risks for people and society. Against this backdrop, policymakers, researchers and practitioners are increasingly interested in public participation in AI: methods that enable members of the public to be involved and have their ideas, beliefs, and values integrated into the design and development process of AI systems [6, 9, 57]. There are two main reasons for this interest: the first is the perceived success of public participation and engagement methodologies in other fields: participatory approaches are used to address issues where there is impact on the public such as in international development [22], environmental justice [42] and in democratic institutions [35]. Increased interest in public participation in AI reflects a broader recognition of AI's implications in the wider world. The second is the, by now, well-documented potential for AI systems to cause harm, such as causing discriminatory impacts on different members of society [4, 18], especially those from marginalised or disadvantaged backgrounds [36, 54]. Proponents of participation cite these methods as a way to create external scrutiny and accountability for these systems [52, 61], and argue 'more or better' participation in AI [46] may partly remedy potential harms [13, 55] and produce more 'socially good' outcomes [13]. Despite this growing interest, it is important to bear in mind that public participation is not a panacea for the harms that AI systems can raise, nor independently capable of deriving societal benefits of emerging technologies. Existing research around 'participation washing' highlights the potential pitfalls and extractive practices of these methods [41, 77].

A review of the literature at the interface between 'participation' and 'AI' reveals that, to date, there is very limited research exploring the role of public participation in commercial AI labs. There is also lingering conceptual confusion about what 'participation' in AI means and what kinds of approaches should be adopted [8, 28], likely hindering wider adoption of these methods. Given that a significant proportion of AI development is undertaken in industry, there is a pressing need to understand how participation is, or could be, embedded in companies driving important developments in AI products and research. This need is all the more urgent in the context of the latest 'AI spring': the advent of novel general purpose and generative AI technologies, which may impact people at greater scale and in more unpredictable ways than traditional 'narrow' AI systems. Tech industry leaders have made calls for more 'public input' into systems like ChatGPT and GPT-4 to ensure these systems are aligned with societal needs [64]. There have also been calls from industry leaders to 'democratise AI', a term that can have different or even conflicting meanings, such as increasing access to these systems or sharing governance of these systems [75]. These developments have intensified the debate about what public participation in AI means.

This paper explores which public participation approaches are being used or considered by tech companies, how they understand the value of these methods,

1

what barriers they face in using these approaches, and what impact public participation has on the company and on participants. Using a literature review of public participation in AI and 12 semi-structured interviews – nine with practitioners working at major AI-focused tech firms, three with non-industry professionals with a stake in the ongoing direction of 'participatory AI' – conducted in the autumn of 2022, this paper seeks to answer three research questions:

**RQ1:** How do commercial AI labs understand public participation in the development of their products and research?

**RQ2:** What approaches to public participation do commercial AI labs adopt?

**RQ3:** What obstacles/challenges do labs face when implementing these approaches?

The contribution of this paper is twofold: novel empirical research reporting perspectives towards and past projects on public participation in commercial AI, and analysis on a current gap in the literature on 'participatory AI', finding that effective uses of participatory methods require a clear understanding of the context in which an AI system will be used.

## 2 Methodology

Our findings emerge from two research verticals: a literature review and semi-structured expert interviews.

### 2.1 Literature review

We surveyed relevant literature on AI ethics and participation, the wider human-computer interaction (HCI), computer supported cooperative work (CSCW) and value-sensitive design (VSD) literature for scholarship on embedding participation in non-AI/ML technologies. We also drew on wider literature focused on the intersections of participation and democracy, for example, including deliberative democracy and sociology. We manually sourced literature from ACM and *arXiv* repositories, using a combination of keyword searches: 'public participation in AI', 'participatory AI', 'participatory design in AI' and 'public engagement', as well as terms and concepts likely to yield discussion of similar/adjacent theoretical grounding including 'social choice' 'and 'democratising AI'. We also used a 'snowball method' to identify additional papers from reference lists.

### 2.2 Expert interviews

We conducted 12 semi-structured interviews in this research. The interviews were led by the lead author, with support and contributions from the second and fourth authors. We interviewed nine practitioners working in large, medium and start-up commercial AI labs developing both products and research, who may be involved in planning or implementation of public engagement /

participation projects or be expected to carry forward findings of public participation projects into research and or product development. For additional background, we also interviewed three subject-matter experts across participatory design, participatory AI and public engagement methods, and with knowledge of tech industry practice. One of these three experts is employed by a technology-focused non-profit, two are currently employed by academic institutions; one of these two had recent previous employment in a commercial lab. All three have authored papers pertaining to participation in AI. See Table 1 for participant IDs. Our interview questions were split into four sections. We asked participants:

1. How they understand public participation;

2. What they think public participation in AI is for;

3. What methods or approaches they have used in their work, or seen in use across the sector, and;

4. Details of their role, their organisation's work culture, resources, and its propensity to fund or conduct participatory work

#### Table 1: Participant organisation and ID

| Organisation | Participant ID |
|---|---|
| Start-up providing open source machine learning | P1 |
| Large company developing both products and research | P2 |
| Large company developing both products and research | P3 |
| Large company developing both products and research | P4 |
| Start-up developing research | P5 |
| Start-up providing open source machine learning | P6 |
| Company developing research | P7 |
| Tech-focused non-profit organisation | P8 |
| Academic institution | P9 |
| Academic institution | P10 |
| Start-up developing products (pre-market) | P11 |
| Company developing research | P12 |

Participants were recruited either directly (selected based on previous demonstrable interest in 'participatory AI', 'responsible AI' or similar fields, and/or were part of the authors' existing industry networks) and through snowball recruitment from recommendations from interviewees. Interviews lasted 60 minutes and took place virtually, using video conferencing software from September 2022 to January 2023, and were transcribed using a speech-to-text transcription software service. Three interviewees did not consent for their interview quotes to be used in this paper. Since all participants were in continuous employment at the time of participation, they were not offered additional payment for their time.

## 2.3 Data analysis

Interview data was analysed using a constructivist qualitative thematic analysis that draws heavily on a 'theoretically flexible' approach set out by Braun and Clarke (2006), that specialises in understanding and reporting repeated patterns, particularly in terms of institutional/organisational behaviours [16]. Using a constructivist epistemology allowed us to approach the data with an understanding that meaning and experience are socially (re)produced [19]. Following this paradigm, we coded our data and constructed our themes according to a 'latent classification' approach [16] surfacing implied beliefs.

The interviews were coded by the lead author using data analysis software. We chose not to set prescriptive benchmarks around prevalence of codes, or whether codes directly related to the RQs. After an initial batch of 71 codes generated, a re-coding process resulted in 56: some codes were felt to be too broad, in other cases, two substantively similar codes were merged (e.g. 'building rapport' to 'relationship building'), and antonyms such as 'inclusion' and 'exclusion' were felt to be usefully interpreted dialectically and coded as single entities.

From these 56 codes, reproduced across Tables 2 and 3, we identified six main themes that corresponded to different research questions:

1. Internal factors

2. Commercial factors

3. Field-level factors

4. Societal and moral factors

5. Purpose of participation

6. Participatory approaches

From the data, we surfaced many different operational considerations and personal values/beliefs that practitioners suggested are (or might be) impactful for the adoption of public participation. Factors were reported to emanate from the level of the firm ('Internal'), or externally ('Field-level'), and pertained to business mission ('Commercial') or relationship to people and society ('Societal and moral'). These are categorised as 'factors' over the more directional e.g. 'blockers' or 'drivers' to avoid setting up a simplistic binary for phenomena not experienced by all participants universally. Some codes appear in different themes, highlighting the porous boundaries between these themes. Theme 5 and Theme 6 concern methods and approaches for, and purpose of, participation, and therefore correspond explicitly with RQ1 and RQ2 of our study.

**Table 2: Themes and codes constructed from factors relevant to the adoption of public participation in commercial AI (as reported by interviewees)**

| Themes | Codes |
| --- | --- |
| Internal factors | Buy-in for public participation<br>Compensating participants<br>Internal expertise<br>Remit: AI product or AI research<br>Responsibility for public participation<br>Scale and scope of public participation<br>Types of 'public'<br>Capacity building |
| Commercial factors | Profit motive<br>PR, optics, reputation<br>Transparency |
| Field-level factors | Capacity building<br>Intermediaries<br>Lack of industry-specific methods or training on public participation<br>PR, optics, reputation<br>Regulation<br>Responsibility for public participation |
| Societal and moral factors | Extractive practice<br>Good intent, social good<br>Harms, discrimination<br>(In)justice, (in)equality<br>Inclusion, exclusion<br>Power<br>Society building<br>Trustworthiness |
| Purpose of participation | Democratising AI<br>Good intent, social good<br>Good business<br>Widening inclusion<br>Embedding lived experience<br>Intrinsic value of participation<br>Public participation as a form of accountability<br>Relationship building<br>Soliciting input / knowledge transfer<br>Trust building |

**Table 2 cont. - Themes and codes constructed from factors relevant to the adoption of public participation in commercial AI (as reported by interviewees).**

| Themes | Codes |
| --- | --- |
| Participatory approaches | Citizens' jury |
| | Crowdsourcing |
| | Co-design |
| | Community training in AI |
| | Community-based approaches |
| | Community-based Systems Dynamics framework |
| | Consultation |
| | Cooperatives |
| | Deliberative approaches |
| | Diverse Voices method |
| | Fairness checklist |
| | Governance tools e.g. audits, impact assessments, other policy mechanisms |
| | Open source |
| | Participatory design |
| | Request for comment |
| | Speculative design/anticipatory futures |
| | Surveys |
| | User research/user testing |
| | Workshops/convenings |

## 2.4 Positionality statement

At the time of research, all the authors were employed by an independent research institute that conducts evidence-based research on data and AI in policy and practice, with a core organisational belief that benefits of data and AI must be justly and equitably distributed, and must enhance individual and social well-being. As part of the organisational remit, the institute collaborates with technology companies in a research capacity, i.e. using industry as a site of study. It does not accept funding from technology companies. The authors live and reside in the UK, and two of the four authors are British, one is British and Irish and one is American. We adopt a sociotechnical conception of AI, understanding that the technical elements of AI – machine learning, neural networks, etc – are inherently interrelated with social, political and cultural factors, principles and motivations (see for example Mohamed et al.) [59].

# 3 Literature review

## 3.1 Public participation in theory and practice

Broadly in the literature, public participation refers to approaches or activities that engage or involve members of the public, incorporating perspectives and experience into a project or intervention. Participatory approaches are routinely adopted in a number of areas, environmental decision-making [43, 50], health and care [62, 72] and in democratic institutions [6, 11]. For example, feedback sessions in health and social care incorporate patient views and lived experience to inform ongoing service delivery (described as 'patient and public involvement' (PPI) in the UK) [10] and consultations in policy mechanisms such as environmental impact assessments foster democratic debate and broaden decision-making powers [43].

In technology design contexts, participatory approaches stem from the fields of human-computer interaction (HCI) [53], user-centred design [57] and the theory and application of participatory design (PD) methods [67]. These fields offer critical examination of how design might be crafted in tandem with [45], instead of on behalf of, different publics in order to incorporate their needs and values [1, 15, 44, 72, 76]. In deliberative democratic theory, it is argued public participation appeals to democratic ideals of legitimacy [78] and accountability [14] as well as to enhance political autonomy [44]. The tradition of deliberative participation – the involvement of the public with a view to fostering deliberative debate and engagement – is evident in participatory design, which offers participants 'seats at the table' [68], emulates democratic decision-making [28], adopts consideration of social and political contexts [1] and embraces co-production [48]. Participation is also often read as an intrinsic value in and of itself [37]: like similar concepts such as 'inclusion' or 'collaboration', it is often understood in the literature as indicative of a 'moral good' [46], of 'flourishing social ties'[17] and so on. However, within the literature, there is little agreement about who constitutes the 'public'. In politics and policy domains, the 'public' may refer to 'citizens', 'labelling data people' or 'laypersons' [42] while, in technology contexts, it may refer to current or future 'end users' [65]. More recent literature around participation in AI adopts a broader definition that includes all people affected by the use of an AI system, particularly individuals and groups for whom AI risks exacerbating inequity, injustice and marginalisation [70]. This raises the question of how commercial AI labs define 'public' in any public participation activities, particularly when their technologies may impact multiple publics in multiple areas or regions.

The form of public participation can vary, reflected in the various typologies produced by political scholars and practitioners [25, 51]. The first of these is Sherry Arnstein's Ladder of Citizen Participation [2], a widely referenced framework for forms of participation, originally intended to outline different degrees of participatory approaches in public planning. Arnstein's eight rungs range from forms of non-participation ('manipulation'), one-way dialogic methods (such as public request for comment [56]), involvement by consultation and partnership in the middle rungs, and finally 'citizen control' at the top rung (see Figure 1). Arnstein is critical of approaches at the bottom of the ladder, branding them tokenistic and inadequate in shifting the axis of power and therefore not paramount to meaningful participation [2, 8].
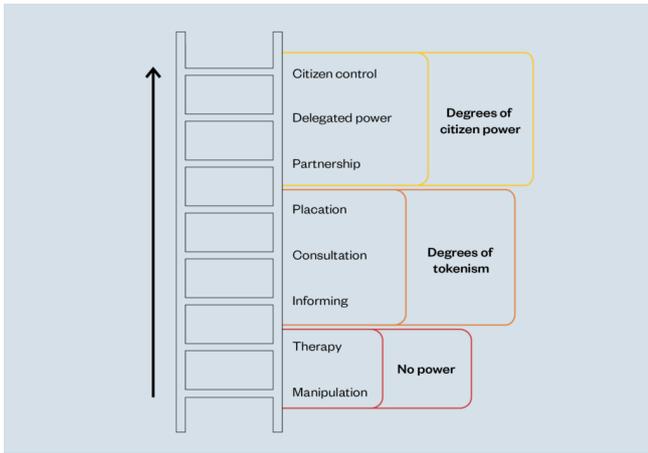
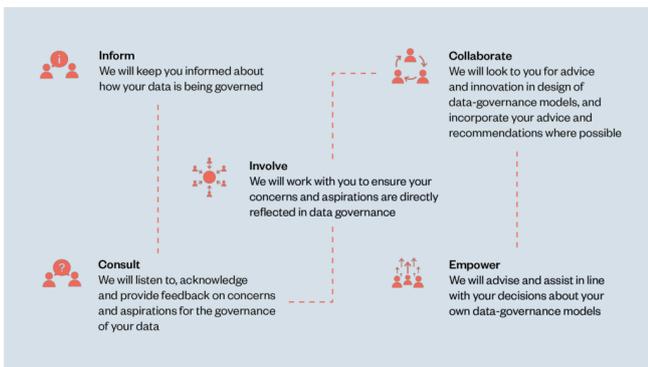**Figure 1: Arnstein's 'Ladder of Citizen Participation' [66] [2]**



**Figure 2: Framework for Participatory Data Stewardship [66]**

Patel et al. [66] draw on Arnstein's ladder and a more recent 'spectrum of participation' [51] to describe practical mechanisms of participation in the stewardship of data and consequently the design of data-driven systems, including AI (see Figure 2). Their analysis creates a link between Arnstein's political lens on participation and participation in sociotechnical contexts by describing five levels of participation and examples of what practical mechanisms may exist for each, drawn from real-world case studies. These five levels include:

1. Informing people about how data about them is used, such as through the publication of model cards;

2. Consulting people to understand their needs and concerns in relation to data use, such as through user experience research or consumer surveys;

3. Involving people in the governance of data, such as through public deliberation or lived experience panels;

4. Collaborating with people in the design of data governance structures and the technologies they relate to, such as through novel institutional structures like 'data trusts', and;

5. Empowering people to make decisions about datasets and technologies built with them, such as through citizen-led governance boards.

Though indirectly linked to AI, these taxonomies help us makes sense of the public participation approaches commercial AI labs may be using and contribute theoretical foundational frameworks for exploring participation in AI design.

## 3.2 Public participation in AI technology development

As Dove et al. note, AI is neither 'arcane nor obscure' [34]: discursive debate around participation in AI should not be isolated from debates around participation more generally. Cooper et al. argue the AI design and development pipeline of AI technologies is diffuse and therefore typically 'participatory', combining multiple iterative activities and the input of multiple actors [24] across 'algorithmic supply chains' [23]. However, as with participation adopted in other domains, there are varying possible degrees of participation in AI. Two existing typologies are instructive for classifying the different modes of participation in AI: Sloane et al.'s typology of participation: *as work, as consultation* and *as justice* [77], and Birhane et al.'s exploration of the three instrumental categories of participation: *for algorithmic performance improvement*; *for process improvement* and *for collective exploration* [8]. These typologies provide a sense of some of the goals of public participation in AI, and where participatory approaches can fit in AI development or research.

There is an emerging literature on participatory approaches to AI development, which identify a few kinds of 'participatory' activities that involve assembling a mixed group of stakeholders to consult or assess an AI system. The literature on participatory development highlights a few activities that are seen as 'participatory'. These include crowdsourcing [31, 81] (such as crowdsourcing possible impacts of ADM systems [5] or labelling data [65]), participatory dataset documentation [80], creating 'red teams' to test or evaluate a model [40], bug bounties [63] or engaging members of the public to elicit preferences for algorithmic design decisions [21, 70]. Such forms of participation very often prioritise a higher total number of participants over length or depth of participant involvement [5]. For example, participatory development of ML datasets [30], requiring higher degrees of input from a higher number of stakeholders might be classified as Sloane's 'participation as work', where methods that foster deliberation around values and experience [32], might fall under Birhane et al.'s heading of 'collective exploration'. Other scholarship argues that participatory approaches in AI could be instrumentalised to advance ambitious societal-level goals such as fairness, inclusion [39, 79], justice [26, 74], accountability [12] and democratic values [38], which could be characterised as Sloane et al.'s 'participation as justice' [77]. Birhane et al. offer three case studies of a participatory approach to AI development, instances where participation is sought to improve the function of large language models for

African and Te Reo Māori languages, annotate datasets and improve dataset documentation [8]. The authors suggest community inclusion in such projects might advance goals such as equity and justice, but acknowledge that participation in these kinds of projects may amount to products built that actually harm the communities included. Another proposed method for participation in AI development is Martin, Jr. et al.'s Community Based System Dynamics (CBSD) method, a mechanism that seeks to 'engage and centre perspectives of marginalized and vulnerable communities' for the purposes of model refinement [58], however only offering cursory detail on the methodological components required to achieve this goal.

There are concerns of 'participation washing' [77] across participation literature, also highlighted in application to AI. Hossain and Ahmed note that, to date, participation in design or development of AI has been overly modest and inconsequential, prescribing only narrow technological solutions as opposed to lasting community or societal change [49], following the general mode of critique from the participation literature [22, 42]. Sloane et al. argue that participatory approaches that claim to value diverse expertise and express a commitment to recentring marginalised communities, but in practice function as (often unrecognised) labour, risk paying lip service to the pro-social ends of participation while exploiting disadvantaged groups [27, 77]. There are also dangers, as noted by Lloyd et al., that with a focus on engaging technology 'users' (in participatory projects), users become a stand-in merely for 'consumers', narrowing focus away from broader segments of society that might be affected by AI, with a risk of exacerbating existing harms to these groups [57]. In instances where a wider focal point is adopted to target 'non-users' of technologies, often under the objective of 'democratising AI' [29], the outcome may not be equivalent to entrenching participatory or democratic structures [73] but may simply indicate intent to 'widening access' to technology use or development [75].

### 3.3 Public participation in commercial AI

Over the past decade, many large technology companies have established or acquired their own dedicated AI labs for developing research and products: for example, the AI research company DeepMind was acquired by Google in 2014 and is now a subsidiary of Google's parent company Alphabet. Google itself has invested in entire AI research wings like Google Brain, and has integrated AI research into its products. There are also a number of smaller, independent companies developing AI that have made significant research and product developments, such as OpenAI and their ChatGPT model and interface. Commercial AI labs are widely considered to be at the forefront of current AI development and research [47].

Many AI labs have teams that are specialised in ethics issues (Microsoft's Office for Responsible AI, Google DeepMind's Ethics and Society team), including a remit for activities such as public participation. Though debates around ethics, fairness and accountability have gained considerable traction in recent years, it is still challenging terrain: Moss and Metcalf point to a habitual inability among firms to to specifically designate which team (members) have the responsibility for embedding ethics [60], as well as an ineptitude toward institutionally buttressing their role(s), creating pinch points and barriers to the effective implementation of AI ethics initiatives. Practitioners struggle with what Rakova et al. identify as a demanding interplay between 'organizational structures and algorithmic responsibility efforts' [69]. Other scholars have criticised tech companies have for 'ethics washing' behaviours, [7] including the use of internal ethics initiatives as a form of social capital that justifies deregulation of their industry in favour of self regulation.

Despite the sheer quantity of industry-led AI/ML research, most scholarship on participation in AI to date has emanated from academia or civil society: there is scant publicly available evidence of what kinds of participatory methods or projects are put into use in commercial AI labs. What literature does exist on public participation approaches in industry is authored by individuals working in commercial AI ethics teams [8, 58], and the limited examples we have of participatory efforts are also led by ethics teams in these companies. Examples include the Royal Society of Arts (RSA) and Google DeepMind's Forum for Ethical AI project, involving a citizens' jury with members of the public to offer space for deliberation on algorithmic decision-making [71], and Behavioural Insights (BIT)'s blog on a recent partnership with Meta constructing citizens' assemblies for members of the public to deliberate on climate misinformation [10]. The lack of public examples of AI labs using participatory methods raises questions about the real extent of their use.

## 4 Interview Findings

Based on our review of the literature, we asked our interview subjects how commercial AI labs understand participatory AI approaches and the obstacles they have faced implementing these practices in the development of AI systems:

1. Within commercial AI labs, public participation is viewed as serving societally 'good' ends, but may also have a strong business purpose;

2. Public participation in AI industry lacks clear and shared understanding of practices. Participants did not identify many participatory methods they use, but rather tended to list methods they had heard of;

3. Public participation in AI labs faces various obstacles: resource-intensity, atomisation, exploitation risk and misaligned incentives;

4. Public participation in AI labs is complicated by products or research that lack a clear context.

## 4.1 Within commercial AI labs, public participation is viewed as serving societally 'good' ends, but may also be good for business

*"We do a lot of AI for social good projects at* [large company]. *But I'm always wondering why we need the qualifier of AI for social good."* [P3]

Interviewees, including the practitioners working on 'participatory AI' and adjacent topics, view participation and participatory approaches positively, with several associating these practices with 'doing good in the world', an indication of company legitimacy or as a commitment to accountability. Another participant described the pull to embed participatory approaches as an 'obligation', to ensure the company are achieving societally beneficial outcomes with their technologies:

*"We, as a corporation, building or researching a technology that has the potential to solve problems for people, have an obligation to engage folks from various backgrounds to help us understand the different problems they face."* [P7]

Some interviewees report viewing public participation in the labs through the lens of profitability or business mission: *"It should be for good business, right? Engaging with people should help you build a product that addresses their wants and needs better which in turn, makes your company more profitable."* [P2]

This view more closely follows the argument that increasing participation in corporate tech contexts presents an opportunity to increase access to technology: unsurprisingly, if your goal is to build better tech, then making it work better for more people is an attractive prospect. However, other participants expressed frustration that this would be likely to be the only logic that would wash with corporate shareholders (who, as one interviewee suggested, would not find any reason to complain if public participation was not conducted at all). In larger companies, interviewees noted challenges of explaining the value and role that participation can play to others in the firm. Those using these methods were trying to resolve concerns around, for example, bias and fairness, but often found that they had to reframe these objectives from the perspectives of how these methods could provide an increasing return on revenue. One interview noted concerns of a performativity around labels such as 'responsible AI':

*"There's concern about being exploitative in using that knowledge to do this sort of marketing veneer of responsible AI, then we're still just going to make money on everything."* [P3]

## 4.2 Public participation in commercial AI lacks a clear and shared understanding of practices

*"We take that there are many different approaches to public participation* [at the company]. *Some are more kind of focused on participatory annotation of data and co-production of AI systems. I think my work is more focused on vision setting for the future of AI."* [P7]

Our research corroborates findings from the literature of an enduring lack of consensus around participatory approaches in practice [8, 26]. Interviewees were asked "what approaches to participation have you used in your work or practice?" Some interviewees were able to talk about approaches they'd personally used for certain research/development projects, but more usually, would recall (often cursory) detail about specific projects or ideas in either their organisation or across the sector, rather than any direct experience. Overall, interviewees cited 1 different methods they were familiar with or had used – see Table 3.

**Table 3: Participatory approaches in commercial AI (as reported by interviewees) mapped onto Arnstein's 'Ladder of Citizen Participation'**

| Arnstein's ladder | Participatory approaches |
|---|---|
| Degrees of citizen power | Cooperatives<br>Citizens' jury<br>Community-based approaches<br>Deliberative approaches<br>Participatory design<br>Speculative design/ anticipatory futures<br>Governance tools e.g. audits, impact assessments, other policy mechanisms |
| Degrees of tokenism | Co-design<br>Community training in AI<br>Community-based Systems Dynamics framework<br>Crowdsourcing<br>UX/user testing<br>Open source<br>Diverse Voices method<br>Workshops/convenings<br>Consultation |
| Non participation | Surveys<br>Request for comment |

The method interviewees cited most often was a form of consultation with people outside the company, generally domain experts rather than members of the public, usually to solicit feedback on the design or usability of products. Most interviewees recognised that participation could have multiple dimensions, with a few specifically using the word 'spectrum'. Two interviewees suggested that open sourcing machine learning models, as a kind of mass participation predicated on widespread involvement, might constitute a participatory approach.

Despite overall understanding and knowledge of types of approaches that could be used in AI development, the important accompanying finding is that most interviewees did not feel fully equipped to report on their organisation's activity in the area of 'participatory AI'. While we cannot rule out that commercial AI labs are using participatory methods that we are unaware of, these

findings suggest that, at best, interviewees did not feel comfortable discussing specific examples of these methods with us, or had no awareness of these methods being used in their companies – and at worst, that such methods are not being used at all. Given that most interviewees self-selected to participate on the basis of their familiarity with public participation in AI (see 'Methodology'), it would appear that the most likely scenario is that there is little use of participatory approaches to AI in industry.

## 4.3 Public participation in commercial AI labs faces various obstacles: resource intensity, atomisation, exploitation risk and misaligned incentives

### 4.3.1 Embedding participation is expensive and resource-intensive

*"If you want actual participation, you actually have to invest before you need something from people."* [P2]

As reported elsewhere in the literature [28], practitioners we spoke with struggled to embed participation in their companies. The accordant time and costs, and the difficulty in quantifying the work, are at present seen as too great to inspire action (and therefore outweighing any motivations for 'social good'). One interviewee put forward interest in conducting further participatory work, but felt that other research and development pursuits, like ensuring 'truthfulness' of large language models, would be a higher priority. Many interviewees put forward a need for capacity building in this space, stating that, at present, practitioners are not equipped to conduct public participation, as many do not come from a social science background or have not undertaken work with community groups, and therefore lack the requisite skills and experience to undertake long-term engagements with members of the public.

### 4.3.2 Participation in the AI industry is 'atomised'

Interviewees often expressed there was not a clear understanding within AI companies of who has the responsibility for leading participatory projects or embedding a 'culture of participation' in which all members of a product team have a shared understanding of the value and uses of these methods. One interviewee suggested that spearheading the adoption of public participation in AI labs puts you at odds with the direction of travel of the rest of the company, effectively creating misaligned incentives, with public participation work not rewarded or recognised within the organisation. In the cases our interviewees mentioned, participation generally arose emergently, responding to specific design or development knots (particularly in the 'agile development' [28] of product lifecycles). One interviewee pointed to burnout and a lack of bandwidth among tech workers, preventing individual practitioners from connecting with other individuals or teams who had taken on participatory work in the past.

### 4.3.3 There is concern and care around exploitation and 'participation washing'

Many interviewees report that they are paying attention to social, societal and moral questions when considering how to adopt public participation approaches in their practice. Frequently cited considerations include concern about extractive behaviour and practice, whether or not 'inclusion' is always a commendable value, and questions of power, justice and societal impacts. Two interviewees specifically cited the term 'participation washing' [77] when sharing thoughts on potential obstacles to embedding participation, which may indicate that this is a concern that has become more routinely observed in these companies. Most interviewees reported feeling great responsibility for non-tokenistic participation and being attuned to power and privilege, especially in capacity as a tech worker. While these interviewees demonstrate a motivation for wanting to adopt meaningful participation that confers decision-making power for participants, for many, it did not translate into 'better' participation (often because owing to the other obstacles we set out in this paper, they felt they could not do a deeper level of engagement justice). Some interview subjects highlighted the tension between the business needs of a commercial lab and the mode of participation in certain projects. While one interviewee reported satisfactory levels of funding and support received by their company, this puts undue pressure on wanting to achieve the 'desired' outcomes from participatory work, recalling a project where they were told to *"go back and get a different answer* [from participants]" [P3]. Other interviewees described concerns of exploitation of participants from marginalised or underrepresented communities in their work:

[recalling previous public participation in the company]*"It gets to the point where it's like 'Oh, yeah, we talked to some Black people. And they said it's fine.' And we're being fair! We're being responsible!"* [P3]

Practitioners report grappling with values such as societal justice and the relation to their work: some discussion across different interviews took place on whether 'inclusion' in AI could advance justice or address power asymmetries. Most interviewees were firm on the importance of adopting focus on communities that have historically been excluded from technology development conversations. For some companies, lowering the barrier of participation/inclusion in AI was deemed a priority, usually in the context of enabling different groups of people to design or use machine learning tools. Moreover, some interviewees situated the role of participation into the broader societal context: one participant argued the role of participation is interrelated to broader questions of political representation and governance: *"That's the realm of the political, setting up the terms under which we all live together. And increasingly, technology, technology systems have encroached so thoroughly on that, that we're having to rethink all of these extremely old*

*questions about how can people self-determine the conditions under which they live in a technology space?"* [P10]

Concern and care over extractive practice and exploitation was reported to closely correlate with the type of 'public' chosen to take part in participatory projects: two interviewees revealed that it is often subject-matter experts that are assembled in place of 'laypeople', suggesting that technical expertise is more often sought out by companies than lived experience. This echoes concerns in the literature around which publics are participating, a particular concern for public participation in AI given the potential for AI systems to impact communities across the globe at great scale and magnitude.

### 4.3.4 Commercial AI labs are not incentivised to be transparent or share their experiences using participatory approaches

Even where participatory approaches are tested and trialled, interviewees described a lack of incentive to report publicly about the work and any potential learnings. One suggested that publishing detail on participatory approaches and specific methodological choices might pose a commercial risk, as it would be sharing information that could be seen as intellectual property.

Some interviewees reported feeling conscious about the reputation of their company, and the ways in which publicising (or not publicising) certain activities could be seen as affecting optics and comprising good or bad 'PR', suggesting that this disincentivises experimenting with public participation.

One interviewee reported feeling as though external scrutiny over practice and public pressure to enact their social responsibilities (where they saw participatory work as situated) did not have much of an effect on the company's direction or bottom line at all: *"If you take all the headlines* [on tech industry practice] *over the last five years, they didn't affect share price, or revenue"* [P3]

This suggests that, for this company, 'the techlash' [3] has not had enormous impact on their practices and would not incentivise publishing details of participatory approaches.

A lack of transparency has effects at the industry level. Institutional theory holds that companies in the sector begin to homogenise when faced with the same set of economic conditions [33], and one interviewee reported that this felt true of tech companies – *"all the AI companies just look at each other "*[P3], suggesting a 'fear of missing out' effect. Coordinated, tech industry-wide effort was often cited by interviewees as being critical for an ecosystem of public participation, particularly around pooling resources to collectively establish or articulate better participatory practices. Most interviewees saw an increased role for some kind of regulation to incentivise public participation, though not without caveat: *"That's a whole other issue of "gaming" regulation. You know, you start this cat and mouse game of: "Here's some regulations". And then companies are thinking, how do we get out of this?"*[P3]

Other actors' contributions to deriving change across the sector was noted by some interviewees, particularly activists. Some suggested looking to other sectors to use as analogues for an AI industry-specific approach. The FDA's medical device pipeline, with its requirement for patient involvement, was offered by two interviewees in this context, as a potential practice that could be adapted to AI research and development.

### 4.3.5 Participation in commercial AI labs is complicated by products or research that lack clear context

As demonstrated above, public participation is costly and resource intensive: companies already lack incentives to conduct it, and where it is conducted, it can be piecemeal. The difficulty of running public participation methods is exacerbated as the generalisability of AI increases.

Three interviewees identified a need to conduct public participation work around more complex, general purpose AI systems where the context in which it could be used to impact the public is less clear, and an additional two were concerned about conducting public participation in the face of rapid development of general purpose AI systems that may present complexities for a non-technical 'public'.

The interviewees we spoke with who belong to or work closely with AI product teams regularly conduct UX/user-research to get feedback on the usability of the proposed product with a narrow group of potential users. Interviewees saw this context as favourable for public input, as potential participants may have a clearer understanding of the impacts of the proposed system: *"Being in a product team can be really focusing, because we have these goals for the conversation. So you can get much clearer feedback from* [participants]" [P2]

One interviewee recalled a project assembling members of the public to discuss potential benefits, harms and use cases of AI models at a high level, but reported that the exercise lacked focus and was not perceived by their company to have useful impact. They suggested that using specific technologies as a steer might enable critical dialogue on possible societal impacts of a technology at a higher level (though did not feel well-equipped to conduct such approaches at present).

Interviewees belonging to research teams, outside strict product deadlines, put forward that they have more flexibility to pursue alternate research or design agendas. For example, practitioners working in research teams had encountered more methods akin to co-design [49, 57] as a result of more agency to set pace and objectives. We find that embedding far-reaching or longer-term public participation projects is seen as particularly complex for general purpose technologies that have many number of downstream applications. One interviewee expressed concern at the pace and spread of recent developments in generative AI further implicating the scope and scale of participation, as well as participant understanding: *"What does it mean to engage people who are affected by, but don't have the knowledge of, state of the art systems, especially as things like DALL-E and DALL-E*

*mini* [now Craiyon] *and Stable Diffusion go viral?*" [P1]

As generative AI and similar technologies continue to proliferate at an astonishing rate, with innumerable downstream uses and a wide user base, several interviewees reported the obligation to conduct some kind of public participation work across a variety of conditions increases, as highlighted by this quote: "*The people that put* [content such as images] *into the public sphere did not know they would be used for this application. How could you know that something you posted in 2007 would be used in a model over a decade later? So the public should have a say.*" [P2]

These findings show that any proposed public participation approach or project must be attuned to the specific context of AI development (product or research). Our findings reveal that it's harder to do public participation when the context in which it would be used or affect the public is less clear (for example, in AI research that is theoretical rather than practical, or with AI systems like generative models that can impact or be used in multiple contexts relevant to a person's life).

# 5 Limitations

## 5.1 Limitations of interview approach

We report the following limitations of our interview approach:

– **Non-representative sample:** Not every major AI lab is represented in this study. In the largest companies, we would have preferred to interview multiple employees from different teams to gain a richer understanding of institutional culture and practice, which is hard to glean from a single interview. Additionally, interviewees in many cases were selected (or self-selected) on the basis of pre-existing interest in ethical/participatory AI etc.

– **Barriers to participation:** We identify two main barriers to participation: interviewee concern around candour, and atomisation of public participation in commercial environments.

Drawing from the research team's prior experience working in industry, and our experiences engaging with industry representatives, we recognised the potential for interviews to surface commercially sensitive IP and or corporate malpractice, resulting in varying degrees of comfort and willingness to interview. Many interviewees may have been reticent to share identifiable details of relevant projects within interview. While we sought to address this limitation by offering interviewees anonymisation of findings and removal of identifiable material, this concern may have persisted. Additionally, as we set out in the Discussion, there is often limited awareness both internally and externally on which individual/team has remit or expertise for public participation, arising in confusion over who would be best placed to participate in this study.

In total, 47 direct personal invitations were sent for this study, in addition to two broadcast messages on two 'responsible tech' Slack boards. 12 directly invited interviewees explicitly declined the offer of participation in this study, we speculate in part owing to some of the barriers set out above, in addition to burnout (which was explicitly cited by a couple of invitees). This resulted in a relatively small sample size of remaining respondents who were available and happy to interview.

## 5.2 Limitations of study

We acknowledge here the recent rounds of tech sector layoffs and the gloomier economic climate beginning to intensify during and shortly following our interview period, and suggest these will have tangible implications for the adoption of participatory approaches (but which are not specifically reported on or studied here). We are employed by a research institute operating in the UK and in Europe, and all interviewees are employed at companies or institutions located in North America and Europe, reflecting the dominant geographies of high-profile AI research labs. We would have preferred to have substantive input from organisations based in the Global Majority represented in this research, though we note, following Chan et al. and others, that mere inclusion is not a conduit to rebalancing North American power domination [20]. Nevertheless, there may be opportunity for future research along these lines.

# 6 Conclusion

In this study, we find that although public participation is recognised as a valuable mechanism to involve public perspectives and enjoys support and interest from this sample of interviewees in commercial AI labs, only limited participatory projects have been explored and implemented to date. Commercial AI labs view public participation as a way to mitigate ethical risks in AI systems and produce more 'societally beneficial' technologies. However, our interviewees report that individuals responsible for implementing participatory approaches in commercial labs do not have a shared understanding of what methods can or should be used and how to use them. While many of the challenges of embedding public participation are not unique to the commercial sector, nor to the context of technology development, there are routinely observed difficulties for public participation in commercial AI: where implemented, participatory approaches in commercial AI labs are informal, atomised and often deprioritised, with limited incentive for companies to publicly declare adoption of participation approaches (even in the context of companies' public commitments to fairness, trustworthiness, and other ethical principles). In some cases, interviewees confirmed concerns from the literature that participation-washing may be occurring.

Consequently, we conclude that factors such as the corporate profit motive and concern around exploitation are at present functioning as significant barriers to the use of participatory methods in AI , rather than drivers or enablers for the uptake of these practices. These con-

cerns for the use of public participation in AI are exacerbated when one considers the growth of general purpose and generative AI systems, which enable a wide range of potential uses of AI systems in different contexts and settings. Successful public participation requires a clear use case for members of the public to understand, raising an innate challenge for the use of these methods for general purpose technologies. It is our intention for this research to function as a springboard: by presenting current conditions and emergent challenges for public participation in commercial AI, we lay foundations for further work and debate.

# 7  Areas for further input

The role of this paper is to provide insight into current challenges in public participation in commercial AI, but this is only one piece in the puzzle in better understanding the logics and conditions of participation in these environments. We acknowledge that possible next steps are manifold, require cooperation from multiple actors, and are unlikely to be 'quick wins'. In light of some of our study limitations, further research on commercial AI public participation is necessary, such as ethnographic research of 'live' participatory projects in labs, to strengthen conclusions on the current lay of the land.

Second, the authors urge industry executives to exercise leadership in this area, namely: connect teams and individuals interested in 'participatory AI' across firms, provide institutional support and funding for further enquiry into participation in AI labs 'in the open' (with learnings made public), and vocally challenge the perceived norm of public participation working in opposition to tech business models. These combined forces may begin to unlock a grander normative vision for what participation in commercial AI should look like.

We join many of our interviewees in their demand for regulators and governments to incentivise this work through appropriate regulatory levers and offer funding and evaluation capacity to kickstart wider adoption of public participation. The authors also recognise and commend the contributions of activists, investigative journalists, researchers and others for their important work in raising awareness of tech industry abuses of power and in advancing algorithmic justice. We call on people affected by uses of AI, activists, civil society and other interest groups to maintain public pressure to advance a stake in the systems and technologies so often built using their data, but decoupled from their values, experiences and vision for technologies and society.

# Acknowledgements

# References

[1] Shana Agid. ""…it's your project, but it's not necessarily your work…": infrastructuring, situatedness, and designing relational practice". In: *Proceedings of the 14th Participatory Design Conference: Full papers - Volume 1.* PDC '16. New York, NY, USA: Association for Computing Machinery, Aug. 15, 2016, pp. 81–90. ISBN: 978-1-4503-4046-5. DOI: 10.1145/2940299.2940317. URL: https://doi.org/10.1145/2940299.2940317 (visited on 04/21/2022).

[2] Sherry R. Arnstein. "A Ladder Of Citizen Participation". In: *Journal of the American Institute of Planners* 35.4 (July 1, 1969). Publisher: Routledge _eprint: https://doi.org/10.1080/01944366908977225, pp. 216–224. ISSN: 0002-8991. DOI: 10.1080/01944366908977225. URL: https://doi.org/10.1080/01944366908977225 (visited on 05/11/2022).

[3] Robert Atkinson et al. *A Policymaker's Guide to the "Techlash"—What It Is and Why It's a Threat to Growth and Progress — ITIF.* 2019. URL: https://itif.org/publications/2019/10/28/policymakers-guide-techlash/ (visited on 05/10/2023).

[4] Imon Banerjee et al. "Reading Race: AI Recognises Patient's Racial Identity In Medical Images". In: *The Lancet Digital Health* 4.6 (June 2022), e406–e414. ISSN: 25897500. DOI: 10.1016/S2589-7500(22)00063-2. arXiv: 2107.10356[cs, eess]. URL: http://arxiv.org/abs/2107.10356 (visited on 01/13/2023).

[5] Julia Barnett and Nicholas Diakopoulos. "Crowdsourcing Impacts: Exploring the Utility of Crowds for Anticipating Societal Impacts of Algorithmic Decision Making". In: *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society.* July 26, 2022, pp. 56–67. DOI: 10.1145/3514094.3534145. arXiv: 2207.09525[cs]. URL: http://arxiv.org/abs/2207.09525 (visited on 09/10/2022).

[6] Aleks Berditchevskaia, Eirini Malliaraki, and Kathy Peach. *Participatory AI for humanitarian innovation: a briefing paper.* London: Nesta, 2021. URL: https://www.nesta.org.uk/report/participatory-ai-humanitarian-innovation-briefing-paper/ (visited on 01/17/2023).

[7] Elettra Bietti. "From ethics washing to ethics bashing: a view on tech ethics from within moral philosophy". In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency.* FAT* '20. New York, NY, USA: Association for Computing Machinery, Jan. 27, 2020, pp. 210–219. ISBN: 978-1-4503-6936-7. DOI: 10.1145/3351095.3372860. URL: https://doi.org/10.1145/3351095.3372860 (visited on 01/10/2023).

[8] Abeba Birhane et al. *Power to the People? Opportunities and Challenges for Participatory AI.* Sept. 15, 2022. DOI: 10.1145/3551624.3555290. arXiv: 2209.07572[cs]. URL: http://arxiv.org/abs/2209.07572 (visited on 09/20/2022).

[9] Abeba Birhane et al. *The Values Encoded in Machine Learning Research.* June 21, 2022. arXiv: 2106.15590[cs]. URL: http://arxiv.org/abs/2106.15590 (visited on 01/07/2023).

[10] BIT. *Deliberative democracy in action.* 2022. URL: https://www.bi.team/blogs/deliberative-democracy-in-action/ (visited on 01/13/2023).

[11] Dan Bloomfield et al. "Deliberation and Inclusion: Vehicles for Increasing Trust in UK Public Governance?" In: *Environment and Planning C: Government and Policy* 19.4 (Aug. 1, 2001). Publisher: SAGE Publications Ltd STM, pp. 501–513. ISSN: 0263-774X. DOI: 10.1068/c6s. URL: https://doi.org/10.1068/c6s (visited on 04/13/2021).

[12] William Boag et al. "Tech Worker Organizing for Power and Accountability". In: *2022 ACM Conference on Fairness, Accountability, and Transparency.* FAccT '22: 2022 ACM Conference on Fairness, Accountability, and Transparency. Seoul Republic of Korea: ACM, June 21, 2022, pp. 452–463. ISBN: 978-1-4503-9352-2. DOI: 10.1145/3531146.3533111. URL: https://dl.acm.org/doi/10.1145/3531146.3533111 (visited on 06/30/2022).

[13] Elizabeth Bondi et al. "Envisioning Communities: A Participatory Approach Towards AI for Social Good". In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. July 21, 2021, pp. 425–436. DOI: 10.1145/3461702.3462612. arXiv: 2105.01774[cs]. URL: http://arxiv.org/abs/2105.01774 (visited on 09/19/2022).

[14] Mark Bovens. "Analysing and Assessing Public Accountability. A Conceptual Framework". In: *European Law Journal* (2007), p. 37. URL: https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1468-0386.2007.00378.x.

[15] Tone Bratteteig and Ina Wagner. "Unpacking the Notion of Participation in Participatory Design". In: *Computer Supported Cooperative Work (CSCW)* 25.6 (Dec. 2016), pp. 425–475. ISSN: 0925-9724, 1573-7551. DOI: 10.1007/s10606-016-9259-4. URL: http://link.springer.com/10.1007/s10606-016-9259-4 (visited on 04/21/2022).

[16] Virginia Braun and Victoria Clarke. "Using thematic analysis in psychology". In: *Qualitative Research in Psychology* 3.2 (Jan. 2006), pp. 77–101. ISSN: 1478-0887, 1478-0895. DOI: 10.1191/1478088706qp063oa. URL: http://www.tandfonline.com/doi/abs/10.1191/1478088706qp063oa (visited on 11/08/2022).

[17] Ellie Brodie, Eddie Cowling, and Nina Nissen. "Understanding participation:" in: *An introduction* (2009), p. 50.

[18] Joy Buolamwini and Timnit Gebru. "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification". In: (2018), p. 15.

[19] Vivien Burr. *An Introduction to Social Constructionism*. Routledge, July 13, 2006. ISBN: 978-0-203-13302-6. DOI: 10.4324/9780203133026. URL: https://www.taylorfrancis.com/books/mono/10.4324/9780203133026/introduction-social-constructionism-vivien-burr (visited on 01/09/2023).

[20] Alan Chan et al. *The Limits of Global Inclusion in AI Development*. Feb. 1, 2021. arXiv: 2102.01265[cs]. URL: http://arxiv.org/abs/2102.01265 (visited on 01/11/2023).

[21] Paul Christiano et al. "Deep reinforcement learning from human preferences". In: (2017). Publisher: arXiv Version Number: 4. DOI: 10.48550/ARXIV.1706.03741. URL: https://arxiv.org/abs/1706.03741 (visited on 05/11/2023).

[22] Frances Cleaver. "Paradoxes of participation: questioning participatory approaches to development". In: *Journal of International Development* 11.4 (1999), pp. 597–612. ISSN: 1099-1328. DOI: 10.1002/(SICI)1099-1328(199906)11:4<597::AID-JID610>3.0.CO;2-Q. URL: https://onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291099-1328%28199906%2911%3A4%3C597%3A%3AAID-JID610%3E3.0.CO%3B2-Q (visited on 01/25/2023).

[23] Jennifer Cobbe, Michael Veale, and Jatinder Singh. *Understanding accountability in algorithmic supply chains*. Rochester, NY, Apr. 7, 2023. URL: https://papers.ssrn.com/abstract=4430778 (visited on 05/04/2023).

[24] A. Feder Cooper et al. "Accountability in an Algorithmic Society: Relationality, Responsibility, and Robustness in Machine Learning". In: *arXiv:2202.05338 [cs]* (May 13, 2022). DOI: 10.1145/3531146.3533150. arXiv: 2202.05338. URL: http://arxiv.org/abs/2202.05338 (visited on 05/19/2022).

[25] Andrea Cornwall. "Unpacking 'Participation': models, meanings and practices". In: *Community Development Journal* 43.3 (July 1, 2008), pp. 269–283. ISSN: 0010-3802. DOI: 10.1093/cdj/bsn010. URL: https://doi.org/10.1093/cdj/bsn010 (visited on 04/11/2023).

[26] Sasha Costanza-Chock. *Design justice: community-led practices to build the worlds we need*. Information policy. Cambridge, MA: The MIT Press, 2020. ISBN: 978-0-262-04345-8.

[27] Nick Couldry and Ulises Ali Mejias. "The decolonial turn in data and technology research: what is at stake and where is it heading?" In: *Information, Communication & Society* 0.0 (Nov. 9, 2021). Publisher: Routledge _eprint: https://doi.org/10.1080/1369118X.2021.1986102, pp. 1–17. ISSN: 1369-118X. DOI: 10.1080/1369118X.2021.1986102. URL: https://doi.org/10.1080/1369118X.2021.1986102 (visited on 01/23/2023).

[28] Fernando Delgado et al. "Stakeholder Participation in AI: Beyond "Add Diverse Stakeholders and Stir"". In: (Nov. 1, 2021). URL: http://arxiv.org/abs/2111.01122 (visited on 04/28/2022).

[29] Paul Dempsey. *Access for all: the democratisation of AI*. Nov. 10, 2021. URL: https://eandt.theiet.org/content/articles/2021/11/access-for-all-the-democratisation-of-ai/ (visited on 02/01/2023).

[30] Emily Denton et al. *Bringing the People Back In: Contesting Benchmark Machine Learning Datasets*. July 14, 2020. arXiv: 2007.07399[cs]. URL: http://arxiv.org/abs/2007.07399 (visited on 01/17/2023).

[31] Mark Diaz et al. "CrowdWorkSheets: Accounting for Individual and Collective Identities Underlying Crowdsourced Dataset Annotation". In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. June 21, 2022, pp. 2342–2351. DOI: 10.1145/3531146.3534647. arXiv: 2206.08931[cs]. URL: http://arxiv.org/abs/2206.08931 (visited on 02/01/2023).

[32] Marc-Antoine Dilhac. *Responsible Artificial Intelligence: a Guide for Deliberation — International observatory on the societal impacts of AI and digital technology*. 2021. URL: https://observatoire-ia.ulaval.ca/en/responsible-artificial-intelligence-a-guide-for-deliberation/ (visited on 05/10/2023).

[33] Paul J. DiMaggio and Walter W. Powell. "The Iron Cage Revisited: Institutional Isomorphism and Collective Rationality in Organizational Fields". In: *American Sociological Review* 48.2 (1983). Publisher: [American Sociological Association, Sage Publications, Inc.], pp. 147–160. ISSN: 0003-1224. DOI: 10.2307/2095101. URL: https://www.jstor.org/stable/2095101 (visited on 01/09/2023).

[34] Graham Dove et al. "UX Design Innovation: Challenges for Working with Machine Learning as a Design Material". In: *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. CHI '17. New York, NY, USA: Association for Computing Machinery, May 2, 2017, pp. 278–288. ISBN: 978-1-4503-4655-9. DOI: 10.1145/3025453.3025739. URL: https://doi.org/10.1145/3025453.3025739 (visited on 01/23/2023).

[35] John S. Dryzek et al. "The crisis of democracy and the science of deliberation". In: *Science* 363.6432 (Mar. 15, 2019). Publisher: American Association for the Advancement of Science, pp. 1144–1146. DOI: 10.1126/science.aaw2694. URL: https://www.science.org/doi/10.1126/science.aaw2694 (visited on 05/10/2023).

[36] Virginia Eubanks. *Automating Inequality: How High-Tech Tools Profile, Police and Punish the Poor*. 2018.

[37] Chad W Flanders. "What Is the Value of Participation?" In: *OKLAHOMA LAW REVIEW* 66 (2013).

[38] Seth Frey, P. M. Krafft, and Brian C. Keegan. ""This Place Does What It Was Built For": Designing Digital Institutions for Participatory Change". In: *Proceedings of the ACM on Human-Computer Interaction* 3 (CSCW Nov. 7, 2019), 32:1–32:31. DOI: 10.1145/3359134. URL: https://doi.org/10.1145/3359134 (visited on 10/14/2022).

[39] William R. Frey et al. "Artificial Intelligence and Inclusion: Formerly Gang-Involved Youth as Domain Experts for Analyzing Unstructured Twitter Data". In: *Social Science Computer Review* 38.1 (Feb. 1, 2020). Publisher: SAGE Publications Inc, pp. 42–56. ISSN: 0894-4393. DOI: 10.1177/0894439318788314. URL: https://doi.org/10.1177/0894439318788314 (visited on 01/11/2023).

[40] Deep Ganguli et al. "Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned". In: (2022).

[41] Michele E. Gilman. *Beyond Window Dressing: Public Participation for Marginalized Communities in the Datafied Society.* Rochester, NY, Nov. 2, 2022. URL: https://papers.ssrn.com/abstract=4266250 (visited on 11/17/2022).

[42] Ludo Glimmerveen, Sierk Ybema, and Henk Nies. "Who Participates in Public Participation? The Exclusionary Effects of Inclusionary Efforts". In: *Administration & Society* 54.4 (Apr. 1, 2022). Publisher: SAGE Publications Inc, pp. 543–574. ISSN: 0095-3997. DOI: 10.1177/00953997211034137. URL: https://doi.org/10.1177/00953997211034137 (visited on 05/12/2022).

[43] Anne N. Glucker et al. "Public participation in environmental impact assessment: why, who and how?" In: *Environmental Impact Assessment Review* 43 (Nov. 2013), pp. 104–111. ISSN: 01959255. DOI: 10.1016/j.eiar.2013.06.003. URL: https://linkinghub.elsevier.com/retrieve/pii/S0195925513000711 (visited on 05/11/2023).

[44] Jürgen Habermas. *Between Facts and Norms: Contributions to a Discourse Theory of Law and Democracy.* Red. by Thomas McCarthy. Trans. by William Rehg. Studies in Contemporary German Social Thought. Cambridge, MA, USA: MIT Press, May 10, 1996. 676 pp. ISBN: 978-0-262-08243-3.

[45] Christina N. Harrington. "The forgotten margins: what is community-based participatory health design telling us?" In: *Interactions* 27.3 (Apr. 17, 2020), pp. 24–29. ISSN: 1072-5520, 1558-3449. DOI: 10.1145/3386381. URL: https://dl.acm.org/doi/10.1145/3386381 (visited on 01/11/2023).

[46] Johannes Himmelreich. "Against "Democratizing AI"". In: *AI & SOCIETY* (Jan. 27, 2022). ISSN: 0951-5666, 1435-5655. DOI: 10.1007/s00146-021-01357-z. URL: https://link.springer.com/10.1007/s00146-021-01357-z (visited on 08/03/2022).

[47] Nathan Benaich {and} Ian Hogarth. *State of AI Report 2022.* 2022. URL: https://www.stateof.ai/ (visited on 02/02/2023).

[48] Kenneth Holstein, Bruce M. McLaren, and Vincent Aleven. "Co-Designing a Real-Time Classroom Orchestration Tool to Support Teacher–AI Complementarity". In: *Journal of Learning Analytics* 6.2 (July 22, 2019). ISSN: 1929-7750. DOI: 10.18608/jla.2019.62.3. URL: https://learning-analytics.info/index.php/JLA/article/view/6336 (visited on 05/11/2023).

[49] Soaad Hossain and Syed Ishtiaque Ahmed. *Towards a New Participatory Approach for Designing Artificial Intelligence and Data-Driven Technologies.* Mar. 30, 2021. DOI: 10.48550/arXiv.2104.04072. arXiv: 2104.04072[cs]. URL: http://arxiv.org/abs/2104.04072 (visited on 08/16/2022).

[50] Stephan Hügel and Anna R. Davies. "Public participation, engagement, and climate change adaptation: A review of the research literature". In: *WIREs Climate Change* 11.4 (2020). _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/wcc.645, e645. ISSN: 1757-7799. DOI: 10.1002/wcc.645. URL: https://onlinelibrary.wiley.com/doi/abs/10.1002/wcc.645 (visited on 05/10/2023).

[51] IAP2. *Core Values, Ethics, Spectrum – The 3 Pillars of Public Participation - International Association for Public Participation.* URL: https://www.iap2.org/page/pillars (visited on 05/04/2022).

[52] Ada Lovelace Institute, AI Now Institute, and Open Government Partnership. *Algorithmic accountability for the public sector.* Ada Lovelace Institute, AI Now Institute, Open Government Partnership, 2021, p. 70. URL: https://www.opengovpartnership.org/documents/algorithmic-accountability-public-sector.

[53] Julie A. Jacko and Andrew Sears, eds. *The human-computer interaction handbook: fundamentals, evolving technologies, and emerging applications.* Human factors and ergonomics. Mahwah, N.J: Lawrence Erlbaum Associates, 2003. 1277 pp. ISBN: 978-0-429-16397-5.

[54] Pratyusha Kalluri. "Don't ask if artificial intelligence is good or fair, ask how it shifts power". In: *Nature* 583.7815 (July 7, 2020), pp. 169–169. DOI: 10.1038/d41586-020-02003-2. URL: https://www.nature.com/articles/d41586-020-02003-2 (visited on 05/16/2022).

[55] Michael Katell et al. "Toward situated interventions for algorithmic equity: lessons from the field". In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency.* FAT* '20. New York, NY, USA: Association for Computing Machinery, Jan. 27, 2020, pp. 45–55. ISBN: 978-1-4503-6936-7. DOI: 10.1145/3351095.3372874. URL: https://dl.acm.org/doi/10.1145/3351095.3372874 (visited on 05/10/2023).

[56] Donald J. Kochan. "The Commenting Power: Agency Accountability through Public Participation". In: *SSRN Electronic Journal* (2017). ISSN: 1556-5068. DOI: 10.2139/ssrn.3006157. URL: https://www.ssrn.com/abstract=3006157 (visited on 03/19/2021).

[57] Alexis Lloyd. *Camera Obscura: Beyond the lens of user-centered design.* Medium. Dec. 21, 2020. URL: https://alexis.medium.com/camera-obscura-beyond-the-lens-of-user-centered-design-631bb4f37594 (visited on 01/11/2023).

[58] Donald Martin Jr. et al. "Participatory Problem Formulation for Fairer Machine Learning Through Community Based System Dynamics". In: *arXiv:2005.07572 [cs, stat]* (May 22, 2020). arXiv: 2005.07572. URL: http://arxiv.org/abs/2005.07572 (visited on 04/13/2022).

[59] Shakir Mohamed, Marie-Therese Png, and William Isaac. "Decolonial AI: Decolonial Theory as Sociotechnical Foresight in Artificial Intelligence". In: *Philosophy & Technology* 33.4 (Dec. 2020), pp. 659–684. ISSN: 2210-5433, 2210-5441. DOI: 10.1007/s13347-020-00405-8. arXiv: 2007.04068[cs, stat]. URL: http://arxiv.org/abs/2007.04068 (visited on 01/31/2023).

[60] Emanuel Moss and Jacob Metcalf. "Ethics Owners: A New Model of Organizational Responsibility in Data-Driven Technology Companies". In: (2020), p. 74. URL: https://datasociety.net/library/ethics-owners/.

[61] Emanuel Moss et al. *Assembling Accountability: Algorithmic Impact Assessment for the Public Interest.* Data & Society, 2021. URL: https://datasociety.net/library/assembling-accountability-algorithmic-impact-assessment-for-the-public-interest/.

[62] Josephine Ocloo et al. "Exploring the theory, barriers and enablers for patient and public involvement across health, social care and patient safety: a protocol for a systematic review of reviews". In: *BMJ Open* 7.10 (Oct. 1, 2017). Publisher: British Medical Journal Publishing Group Section: Health services research, e018426. ISSN: 2044-6055, 2044-6055. DOI: 10.1136/bmjopen-2017-018426. URL: https://bmjopen.bmj.com/content/7/10/e018426 (visited on 11/22/2021).

[63] OpenAI. *Announcing OpenAI's Bug Bounty Program.* 2023. URL: https://openai.com/blog/bug-bounty-program (visited on 05/10/2023).

[64] OpenAI. *How should AI systems behave, and who should decide?* Feb. 17, 2023. URL: https://openai.com/blog/how-should-ai-systems-behave (visited on 04/20/2023).

[65] Junwon Park et al. "AI-Based Request Augmentation to Increase Crowdsourcing Participation". In: *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 7 (Oct. 28, 2019), pp. 115–124. ISSN: 2769-1349. DOI: 10.1609/hcomp.v7i1.5282. URL: https://ojs.aaai.org/index.php/HCOMP/article/view/5282 (visited on 02/01/2023).

[66] Reema Patel et al. *Participatory data stewardship*. Ada Lovelace Institute, 2021. URL: https : / / www . adalovelaceinstitute.org/report/participatory-data-stewardship/ (visited on 01/10/2022).

[67] Jennifer Pierre et al. "Getting Ourselves Together: Data-centered participatory design research & epistemic burden". In: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. CHI '21: CHI Conference on Human Factors in Computing Systems. Yokohama Japan: ACM, May 6, 2021, pp. 1–11. ISBN: 978-1-4503-8096-6. DOI: 10.1145/3411764.3445103. URL: https://dl.acm.org/doi/10.1145/3411764.3445103 (visited on 06/29/2022).

[68] Scott Probiner and Timothy Murphy. *From smart products to smart systems*. 2018. URL: https : / / www2 . deloitte . com / content / www / us / en / insights / focus / cognitive - technologies / participatory - design - artificial - intelligence.html (visited on 01/31/2023).

[69] Bogdana Rakova et al. "Where Responsible AI meets Reality: Practitioner Perspectives on Enablers for shifting Organizational Practices". In: *Proceedings of the ACM on Human-Computer Interaction* 5 (CSCW1 Apr. 13, 2021), pp. 1–23. ISSN: 2573-0142. DOI: 10.1145/3449081. arXiv: 2006.12358[cs]. URL: http://arxiv.org/abs/2006.12358 (visited on 04/11/2023).

[70] Samantha Robertson and Niloufar Salehi. "What If I Don't Like Any Of The Choices? The Limits of Preference Elicitation for Participatory Algorithm Design". In: *arXiv:2007.06718 [cs]* (July 13, 2020). arXiv: 2007.06718. URL: http : / / arxiv . org / abs / 2007 . 06718 (visited on 04/29/2022).

[71] The RSA. *Democratising decisions about technology: a toolkit*. Oct. 24, 2019. URL: https : / / www . thersa . org / reports / democratising - decisions - technology - toolkit (visited on 02/03/2023).

[72] Jill Russell, Nina Fudge, and Trish Greenhalgh. "The impact of public involvement in health research: what are we measuring? Why are we measuring it? Should we stop measuring it?" In: *Research Involvement and Engagement* 6.1 (Dec. 2020), p. 63. ISSN: 2056-7529. DOI: 10.1186/s40900-020-00239-w. URL: https://researchinvolvement.biomedcentral . com / articles / 10 . 1186 / s40900 - 020 - 00239-w (visited on 05/11/2023).

[73] Henrik Skaug Sætra, Harald Borgebund, and Mark Coeckelbergh. "Avoid diluting democracy by algorithms". In: *Nature Machine Intelligence* 4.10 (Sept. 29, 2022), pp. 804–806. ISSN: 2522-5839. DOI: 10.1038/s42256-022-00537-w. URL: https://www.nature.com/articles/s42256-022-00537-w (visited on 05/11/2023).

[74] Kristen M. Scott et al. "Algorithmic Tools in Public Employment Services: Towards a Jobseeker-Centric Perspective". In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. FAccT '22: 2022 ACM Conference on Fairness, Accountability, and Transparency. Seoul Republic of Korea: ACM, June 21, 2022, pp. 2138–2148. ISBN: 978-1-4503-9352-2. DOI: 10.1145/3531146.3534631. URL: https://dl.acm.org/doi/10.1145/3531146.3534631 (visited on 08/02/2022).

[75] Elizabeth Seger et al. *Democratising AI: Multiple Meanings, Goals, and Methods*. Mar. 27, 2023. DOI: 10.48550/arXiv.2303.12642. arXiv: 2303.12642[cs]. URL: http://arxiv.org/abs/2303.12642 (visited on 04/26/2023).

[76] Mona Sloane. "To make AI fair, here's what we must learn to do". In: *Nature* 605.7908 (May 5, 2022), pp. 9–9. ISSN: 0028-0836, 1476-4687. DOI: 10.1038/d41586-022-01202-3. URL: https://www.nature.com/articles/d41586-022-01202-3 (visited on 05/11/2023).

[77] Mona Sloane et al. *Participation is not a Design Fix for Machine Learning*. Aug. 11, 2020. arXiv: 2007.02423[cs]. URL: http : / / arxiv . org / abs / 2007 . 02423 (visited on 05/11/2023).

[78] Stephanie Solomon and Julia Abelson. "Why and When Should We Use Public Deliberation?" In: *Hastings Center Report* 42.2 (2012), pp. 17–20. ISSN: 1552-146X. DOI: 10.1002/hast.27. URL: https://onlinelibrary.wiley.com/doi/abs/10.1002/hast.27 (visited on 05/12/2022).

[79] Jack Stilgoe, Richard Owen, and Phil Macnaghten. "Developing a framework for responsible innovation". In: *Research Policy* 42.9 (Nov. 1, 2013), pp. 1568–1580. ISSN: 0048-7333. DOI: 10.1016/j.respol.2013.05.008. URL: https://www.sciencedirect.com/science/article/pii/S0048733313000930 (visited on 03/22/2021).

[80] Harini Suresh et al. "Towards Intersectional Feminist and Participatory ML: A Case Study in Supporting Feminicide Counterdata Collection". In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. FAccT '22: 2022 ACM Conference on Fairness, Accountability, and Transparency. Seoul Republic of Korea: ACM, June 21, 2022, pp. 667–678. ISBN: 978-1-4503-9352-2. DOI: 10.1145/3531146.3533132. URL: https://dl.acm.org/doi/10.1145/3531146.3533132 (visited on 06/29/2022).

[81] Jennifer Wortman Vaughan. "Making Better Use of the Crowd: How Crowdsourcing Can Advance Machine Learning Research". In: (2017), p. 46. URL: https://jmlr.org/papers/v18/17-234.html.