

# Understanding AI Data Production & Community Impacts Worldwide: A Multivocal Literature Review

ANONYMOUS AUTHOR(S)

Artificial intelligence (AI) depends on data production: the sociotechnical process that transforms human knowledge into computational resources. The connections among AI systems, data practices, and impacts on Indigenous, underrepresented, and underserved communities—though critical—have not been systematically examined. To this end, we conduct a Multivocal Literature Review (MLR) integrating 350 academic and grey-literature sources to analyze how AI systems, data practices, and community impacts intersect. Across five analytic domains—Data Relations, Data Labor, Data Representation, Data Infrastructure, and Data Governance—we distinguish extractive data production mechanisms that centralize control from high-agency pathways in which communities exercise authority. We contribute (1) a multivocal review that positions data production as a site of sociotechnical power rather than a technical prerequisite; (2) implications for HCI research including upstream infrastructure as a design site, provenance-first architectures, and federated data governance supporting community sovereignty; (3) methodological validation of multivocal synthesis for bridging academic critique with community practice; and (4) an open corpus mapping sources across pipeline stages, historical eras, and geographic contexts.

CCS Concepts: • **Do Not Use This Code → Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

Additional Key Words and Phrases: AI, ML pipeline, data production, extractive practices, underserved communities, Indigenous data sovereignty, data collection

## ACM Reference Format:

Anonymous Author(s). 2018. Understanding AI Data Production & Community Impacts Worldwide: A Multivocal Literature Review. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 33 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 Introduction

Contemporary artificial intelligence (AI) systems depend on data. As approaches have advanced over the past three decades, the scale and composition of data needs has transformed: from small expert-curated datasets like MNIST [86], to massive crowdsourced benchmarks such as ImageNet [40], and now to foundation models trained on billions of scraped web documents, images, and interaction traces [131, 144]. The opacity and complexity of the machine learning (ML) pipeline [23], as well as the diversity and amount of human knowledge and labor needed [129, 166], has expanded dramatically.

This trajectory matters because the choices made in gathering and curating data directly shape which communities benefit from AI systems and which communities bear their costs. Data is made, not found—it is produced through a series of choices about what to gather, how to curate it, and under what terms. Decisions made upstream and all along the pipeline can either create or mitigate harms, which fall disproportionately on underserved, underrepresented,

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

and Indigenous communities [7, 14, 80, 90, 103, 126]. Despite growing critical attention to algorithmic harms and dataset bias, data production itself—the complex sociotechnical process through which human knowledge becomes computational input—has rarely been treated as a central object of inquiry within HCI.

Relevant literature is scattered across disciplines and publication ecosystems. Human–computer interaction (HCI) contributes a long-standing body of research on how sociotechnical systems enact power, from canonical postcolonial critiques [68, 157] to recent work analyzing “extractive” dynamics in ICT4D research [47], and epistemic injustice [5, 171]. Review literature in HCI and adjacent fields offers substantial insight into how AI systems affect communities. A scoping review by Shelby et al. [149] maps harms experienced by underserved groups, and a subsequent systematic review by Wang et al. [168] builds on this foundation through a focused synthesis across disability contexts. Case-based analyses examine how AI systems mismatch their contexts of use ([138]. Surveys of AI ethics and NLP bias identify structural gaps and the absence of lived-experience perspectives [17, 19], while dataset genealogies trace how collection and annotation practices embed exclusions [41]. Recent HCI scholarship analyzes how misabstraction cascades through sociotechnical systems [37] and takes stock of social-justice commitments within HCI [28].

This body of work illuminates important dimensions of a three-part relationship: AI systems, data production practices, and community impacts. Yet, most scholarship examines discrete pieces rather than synthesizing across all three. Most reviews also draw primarily on academic publications, capturing the scientific state-of-the-art but leaving less visible documentation of the state-of-practice. These gaps call for synthesis across this three-part relationship. In response, we use a multivocal literature review (MLR), an approach designed to synthesize knowledge that circulates across many publication ecosystems. An MLR integrates white literature (peer-reviewed academic publications) and grey literature (policy reports, organizational materials, community outputs), offering a structured way to work with a wider evidence landscape [53]; see Appendix A.

We treat *data production* as the complex sociotechnical process through which data is defined, gathered, curated, and controlled across model pipelines.<sup>1</sup> This framing supports a move beyond understanding *data collection* as a routine methodological disclosure or neutral technical artifact. Instead, we center the institutional choices, power relations, and consequences that underpin AI development and determine who benefits from or bears its costs. Our approach leverages Critical Computing as a diagnostic lens and Social Justice as a normative orientation. Critical Computing shows how data practices reflect institutional priorities and labor arrangements rather than some objective “ground truth” modeled by engineers, and offers tools for analyzing power in dataset construction and use. Social Justice complements the diagnosis by asking how data work might redistribute agency, benefit, and governance toward the communities whose knowledge and labor support AI systems. Together, the two orientations clarify why upstream data production is a sociotechnical domain of timely concern for HCI and a site for upstream design intervention.

Overall, this work provides a literature review with novel concepts of data production as a sociotechnical site where power is negotiated and encoded, rather than a logistical preliminary to model development. We identify mechanisms where extractive practices and high-agency alternatives diverge across the ML pipeline, with implications for HCI. In summary, we contribute:

- A multivocal literature review of 350 sources across academic and grey literature, resulting in five analytic domains where AI systems, data production, and community impacts intersect;

<sup>1</sup>We share with Miceli & Posada [99] an emphasis on “production” to foreground relations of power and knowledge in data and labor, which echoes the “assemblage” approach of Kitchin et al. [78], also rooted in Foucauldian critique.

- Opportunities for HCI research and practice, including upstream data infrastructure as a design site, provenance-first architectures, federated learning for community sovereignty, and ethics review paradigms that scrutinize data production;
- Methodological validation of multivocal synthesis for sociotechnical inquiry, showing how grey literature captures state-of-practice missing from academic venues;
- An open corpus of 350 sampled sources mapped across pipeline stages, historical eras, and geographic contexts, with structured summaries and documented rationales for relevance to the three-part inquiry, publicly available at <https://github.com/ADC-chi/ai-data-production-landscape>.

## 1.1 Key Terms and Definitions

*1.1.1 Artificial Intelligence.* For clarity, we use “AI” in this paper primarily in its modern ML sense. An ML system learns patterns and rules from training data to create a predictive model [23]. The resulting model must then be evaluated for reliability and generalization using a separate, independent dataset known as test data [44]. We contextualize our discussion of AI within its broader historical trajectory of research and development but focus on the current statistical and data-driven era of AI that facilitates many contemporary extractive regimes [57]. We intend our discussion to be situated not merely as a critique of modern ML but as a reflection on a continuous thread within technological and social history.

*1.1.2 Extractive.* We use “extractive” in this paper to denote high-asymmetry or dispossessive practices, building upon its conventional association with Indigenous marginalization and digital forms of resource appropriation. This definition is designed to encompass diverse historical and contemporary manifestations of power imbalances that result in one party’s advantage at the expense of another’s autonomy or resources [20, 117]. Illustrative examples of such high-asymmetry practices are evident in historical contexts, such as the exploitative labor practices of UK coal mining [35]; the profoundly unethical nature of the Tuskegee syphilis experiment in the United States and the untreated carcinoma study in New Zealand [72, 124]; and contemporary issues like large-scale industrial mining [110] and pervasive AI surveillance [121]. By broadening this definition, our objective is to more accurately encapsulate the systemic character of extraction across various domains. In contrast, we describe as high-agency examples of principles and practices in the literature which prioritize active participation, equitable distribution of power, and community-defined obligations [127, 174]. These examples appear in contexts where communities, practitioners, or institutions negotiate shared authority, shape the terms of data contribution, or establish governance arrangements that align data use with locally grounded priorities.

*1.1.3 Communities and Populations.* We use “underserved” to describe communities lacking adequate infrastructural, institutional, or economic support, and “underrepresented” to indicate groups whose knowledge, languages, or perspectives are numerically absent or devalued in AI research and development [91, 149]. We use the umbrella category “Indigenous,” which “enables historically and geographically separated peoples to recognize each other and their common plight, and to collaborate towards a better future” [135]. We avoid “marginalized” in the adjectival form to emphasize agency and resistance rather than positioning communities as passive victims. We use Global Majority to emphasize that most of the world’s population lies outside Euro-American contexts. Our chosen terms underscore structural asymmetries in power and resource distribution rather than deficits within communities themselves [158].

## 2 Background

### 2.1 Critical Traditions on Extraction and Justice

Foundational works from theoretical, historical, and community traditions establish frameworks for studying power in knowledge production. Theories of epistemic violence and injustice [49, 153] and “situated knowledges” [64] interrogate how knowledge systems encode relations of domination. Historical analyses of colonial resistance show alternative epistemologies and organizing strategies [70].

Black feminist theory articulates intersectional approaches to structural power, from early collective statements [31] to analyses of interlocking systems of oppression [32]. Gender and queer theory establish frameworks for analyzing the production of normativity [24], binary logics [147], and classificatory power [30]. Indigenous studies center community sovereignty and relational ethics [10, 39, 89] and provide frameworks for decolonizing knowledge production in research [152] and AI data practices [21]. Critical data studies crystallize a complementary set of concerns for the digital context, with a focus on datafication, surveillance, and governance [79].

Foundational works attune us to centuries of extractive patterns, resistance, and knowledge-making. They are essential for understanding present and future technological worlds. Here, critical works anchor the conceptual vocabulary of extraction and justice in the context of the global AI data production ecosystem.

### 2.2 Evolution of AI Data Production Practices

Since the advent of machine learning, there has been a constant need for data. Over time, how that data was produced has undergone transformations beyond dataset sizes. These changes include how data is produced and who performs the work [41]. As demands for larger models have intensified, practices have shifted from small, carefully curated corpora, to large datasets assembled through web-scraping and crowdsourced annotations, to massive, automated web-scraped collections supported by industrial-scale annotation.

*Era 1.* Early curated datasets were small, domain-specific, and selected by experts. A canonical example is MNIST, a dataset of handwritten digits drawn from U.S. postal codes [85]. Choices about inclusion and categorization reflected institutional knowledge and disciplinary priorities.

*Era 2.* Large curated datasets expanded scale through crowdsourced annotation of web-scraped content, exemplified by ImageNet [40] and MS COCO [93]. This era accelerated deep learning [38, 58] but shifted labor from domain experts to distributed workers, often in Global Majority regions, as well as automated web-scraping efforts, ultimately emphasizing performance gains over contextual fit.

*Era 3.* Contemporary data production diverges into two parallel approaches. Massive, largely uncured web-scraped corpora such as Common Crawl [9, 33] and C4 [43, 131], LAION [144, 145], Refined Web [125], and ClueWeb22 [116] are assembled through automated scraping at an unprecedented scale. Such production efforts shape and are shaped by competitive foundation model development [16, 156]. Alongside and often in response, smaller, highly curated datasets emerged, produced through participatory methods and community partnerships. Examples include ROOTS [84], Masakhane’s African language collections [112], and Cohere’s multilingual Aya Dataset [151].

Dataset hosting and governance practices have shifted over time as well: from freely downloadable units like MNIST, to single-location storage on cloud services (e.g., AWS, HuggingFace), to URL-indexed collections like LAION that disclaim responsibility for original sources, and to emerging federated “data spaces” designed to support locally owned infrastructures and community governance [66]. Proprietary datasets are closed, and open-source alternatives range

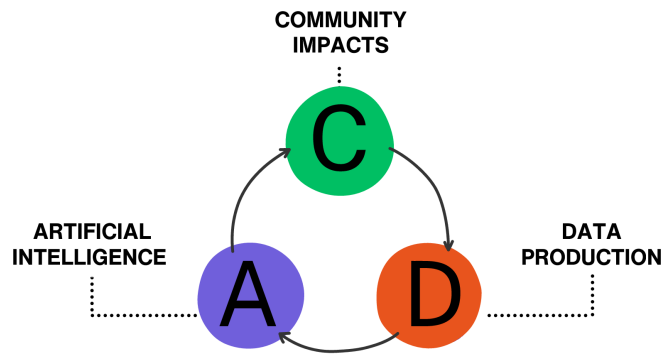


Fig. 1. A/D/C scoping framework for identifying relevant literature across three interrelated domains of AI (A), Data Production (D), and Community Impacts (C).

from massive scrapes to carefully stewarded community collections; each modality comes with distinct risks and obligations [92, 172].

The three eras feature distinct technical capabilities, institutional arrangements, and methodologies that enabled extractive patterns to scale industrially. While both Era 2 and Era 3 are founded in web-scraped datasets, the scale at which Era 3 extracts data is unprecedented. As such, the current era hosts both the most expansive extractive practices and the most developed community-controlled frameworks. The future of AI data production is not determined.

### 3 Methodology

We conducted an MLR to assemble a structured body of evidence on AI data production and its impacts on underserved, underrepresented, and Indigenous communities. We treated academic and grey sources as complementary evidence streams. The research-design diagram and other accompanying materials are available on the companion site: <https://adchi.github.io/ai-data-production-landscape/>.

Because our inquiry has three interrelated parts—AI systems, data production practices, and community impacts—we formalized it through what we call the A/D/C framework (Figure 1). This framework defined the scope of our inquiry, which was ultimately bounded by community impacts (C). All sources engage substantively with community experiences, power dynamics, or consequences in contexts relevant to AI data production. Sources varied in whether they directly addressed AI systems (A) and data production practices (D), or provided foundational understanding that informs interpretation of these dimensions. Section 3.2 provides more detail on how we created the corpus based on this framework.

#### 3.1 Search Strategy

We developed the search strategy by deriving keywords from the A/D/C framework. We searched seven academic databases between August 2024 and January 31 2025: ACM Digital Library, IEEE Xplore, ScienceDirect, Taylor & Francis Online, Wiley Online Library, Springer Link, and Google Scholar. We iteratively developed boolean search strings with AND/OR terms across variants of A/D/C terms using Boolean operators. Titles and abstracts were screened first,

followed by full-text assessment for items meeting initial criteria. Table 1 shows the key and supplementary search terms of our inquiry.

Table 1. Key and Supplementary Search Terms

| Dimension                    | Key Terms   | Supplementary Terms  |
|------------------------------|---|--|
| <b>AI Systems (A)</b>        | Artificial Intelligence, Machine Learning, AI, ML                                 | Large Language Models, LLMs, Computer Vision, Foundation Models, Neural Networks, Automated Systems, Algorithmic Systems, Deep Learning          |
| <b>Data Production (D)</b>   | Data Collection, Data Production, Dataset Creation, Data Curation, Data Practices | Annotation, Labeling, Data Labor, Crowdsourcing, Web Scraping, Data Extraction, Dataset Development, Data Work, Data Gathering, Responsible AI   |
| <b>Community Impacts (C)</b> | Indigenous, Marginalized, Underrepresented, Underserved, Community                | Global Majority, Global South, Data Sovereignty, Linguistic Diversity, Cultural Context, Extraction, Appropriation, Bias, Fairness, Harm, Safety |

We developed four primary query sets: foundational (targeting core data production practices in AI contexts affecting communities), extraction frame (targeting exploitative practices), data labor (focusing on crowdsourcing and platform labor), and alternatives (seeking participatory and community-led approaches). The ACM Digital Library search illustrates our results. Across the four query sets, 1,914 hits yielded 1,201 items screened, 153 meeting initial criteria, and 48 unique sources after full-text review and duplicate removal. Similar strategies applied to the remaining six databases. Database searches contributed 174 sources, representing 50% of the final corpus. See Appendix B for more search details.

For grey literature we used different methods. Following Garousi et al. [53], we used general Google Search and systematically examined organizational ecosystems engaged in AI data work, prioritizing organizational reports, policy documents, and community outputs from established entities. Three complementary methods supplemented database and grey literature searches. Citation snowballing [173] from 20 seed papers tracked forward and backward citations iteratively, contributing 51 sources (15%). Hand-searching of journals including CHI, FAccT, CSCW, *Journal on Responsible Computing*, *ACL*, *Big Data & Society*, and *AI & Society* contributed 21 sources (6%). Iterative gap-filling searches addressed underrepresented regions, concepts, or pipeline stages as the corpus took shape, contributing 31 sources (9%).

Inclusion criteria followed from the A/D/C framework. Sources entered the corpus when they engaged community impacts substantively. Most sources additionally provided direct evidence about data production practices or AI systems. This meant including sources that analyzed AI system behavior, deployment, or evaluation in relation to community outcomes; examined data sourcing, processing, annotation, governance, or infrastructure with implications for affected communities; provided community-governed protocols, sovereignty statements, or governance frameworks; or established theoretical or epistemological foundations addressing power, resistance, extraction, or marginalization in ways essential for interpreting AI data practices and their consequences. We excluded sources that, for example, discussed AI ethics, fairness, or responsible AI at a high level without addressing data practices or community impacts; focused solely on model performance, technical optimization, or algorithmic advances without sociotechnical analysis; reported community-based research unrelated to AI systems or data production; or were non-English (a pragmatic limitation we discuss more below).

Quality assessment varied by source type and followed guidance from Kamei et al. [74]. Academic items underwent venue peer review. Grey-literature items required additional evaluation; we assessed organizational authority, author expertise, community recognition, and the provenance of policy documents, and we interpreted community outputs through alignment with decolonizing methodologies and community endorsement. These criteria track with grey-literature appraisal in multivocal reviews and draw on elements of Garousi et al. [53]’s framework, including stated aims, methodological clarity, contribution, and outlet type, which for this work primarily meant community-governed and sovereignty-oriented materials.

Screening proceeded in two stages: titles and abstracts were reviewed for relevance to the C boundary, followed by full-text assessment. The first author led database and grey literature searches. Two authors independently read full-text articles, prepared summaries, and presented sources in batches to the full team for consensus review. Disagreements on inclusion criteria and relevance to the A/D/C framework were resolved through discussion. This process occurred in two rounds (August–September 2024), with each round reviewing approximately 100 candidate sources. These consensus rounds established shared standards before the two authors completed full screening of the 350-source corpus in February 2025. The final corpus contains 258 academic items (74%) and 92 grey-literature items (26%).

### 3.2 Corpus Creation

We created a datasheet that categorizes each source across multiple dimensions to provide maximum contextualization [37]. Coded categories included bibliographic metadata (author, year, venue, type), A/D/C coverage, pipeline stage, historical era, orientation, geographic focus, and author affiliation. For each source, we additionally recorded a unique rationale for A/D/C coverage and a brief summary to support traceability of sourcing and selection decisions. The same two authors who led source selection conducted categorization using the collaborative consensus approach established during screening.

**A/D/C Coverage.** We categorized each source based on where direct evidence appears across our three-part inquiry: AI systems (A), data production practices (D), and community impacts (C). Every source engages all three dimensions analytically, but sources vary in what they directly support versus what requires interpretive connection. We wrote rationales stating what each source contributes to A, D, and C to clarify where direct evidence appears and where relevance is interpretive and to provide transparent disclosure of our interpretive stance on each source. Tags indicate where direct evidence is present. We do not tag A or D dimensions alone because all sources must engage community impacts (C) to enter the corpus. Our coding produced four categories, described in Table 2.

**Pipeline Stages.** We mapped each source to stages of a simplified AI development pipeline (Figure 2): *Problem Understanding & Formulation* (institutional prioritization, funding decisions, and product conception), *ML System Design and Development* (data selection and enrichment, model architecture choices, and training processes), and *Deployment & Impact* (product testing, launch, and post-deployment effects) [98]. We mapped each source to a pipeline stage and sub-stage to make visible where specific mechanisms arise and how decisions at those points propagate through later phases—what prior work characterizes as cascading effects that compound downstream harms [136, 155]. Sources spanning multiple stages or describing cross-cutting dynamics were tagged accordingly.

**Historical Eras.** We distinguished three eras of data production, per discussion in 2.2: Era 1 (expert-curated datasets, pre-2009), Era 2 (crowdsourced benchmarks, 2009–2017), and Era 3 (web-scraped and foundation models, 2017–present). Multi-era sources were coded accordingly. No sources were coded exclusively as Era 1.

**Orientations.** Each source received a single orientation code reached through team consensus based on the primary analytical purpose the source served in this inquiry. *Extractive* sources provided direct evidence of practices undermining

Table 2. Triangle coverage coding definitions showing how corpus sources engage AI systems (A), data production practices (D), and community impacts (C)

| Code | Description   | Example Sources   |
|------|---|---|
| ADC  | Direct evidence relevant to AI systems, data production, and community impacts            | Garcia et al. [51] on critical refusal as an intervention into extractive data logics and governance; Hall et al. [59] on participatory, community-engaged dataset production; Park et al. [120] on designing accessible infrastructures for collecting AI data from people with disabilities; Rifat et al. [133] on categorization politics and context erasure in annotating faith-based violence data; Lewis et al. [89] on Indigenous protocol-aligned dataset construction and culturally grounded AI applications |
| DC   | Direct evidence for data production and community impacts; AI relevance is interpretive   | Adley et al. [4] on ethical data collection with marginalized groups and power dynamics in practice; Cooper et al. [34] on community-collaborative research models emphasizing shared control and benefit; Hancock et al. [60] on tensions in data sharing and harms within a modern slavery data ecosystem; Taylor and Kukutai [159] on Indigenous Data Sovereignty and metadata governance; Pool [128] on colonial census practices replacing Māori knowledge systems   |
| AD   | Direct evidence for AI systems and data production; community impacts are clearly implied | Bhardwaj et al. [12] on evaluating ML datasets through a data-curation lens and FAIR principles; Koch et al. [82] on dataset reuse and benchmark concentration; Sambasivan et al. [136] on data cascades and hidden labor in high-stakes ML pipelines; Schiff et al. [143] on translating AI principles into practice via participatory, iterative impact assessment; Zhao et al. [177] on fairness-curation challenges faced by dataset curators across organizational and socio-political contexts                    |
| C    | Direct evidence about community impacts only; A and D relevance is interpretive           | Battiste [10] on Indigenous epistemologies and marginalization; Haraway [64] on situated knowledge and partial perspective; Igwe et al. [67] on non-extractive research principles; James [70] on colonial extraction economies and collective resistance in the Haitian Revolution; Shapiro and McNeish [148] on hyper-extractivism and resistance   |

consent, compensation, or community benefit. *High-agency principles* advanced normative frameworks with explicit policy or governance recommendations. *High-agency practices* described operationalized initiatives with concrete implementation details.

**Synthesis.** We synthesized findings through iterative analysis across these dimensions. When multiple sources described similar mechanisms across different contexts, we consolidated these into recurring patterns. Individual sources could exemplify multiple patterns. Patterns were distilled into the five analytic domains described in Section 4. Complete coding definitions with examples appear in Appendix C.

### 3.3 Limitations and Reflexivity

We recognize that subjectivity shapes our interpretations, though transparent documentation and multi-researcher validation helped address this inherent limitation. Connecting multiple disciplinary traditions, historical eras, and global contexts proved challenging, creating “translation needs” across distinct vocabularies and epistemological frameworks. The sourcing strategy privileged networks in Africa and global Indigenous movements, yielding detailed coverage of

# AI Development Pipeline

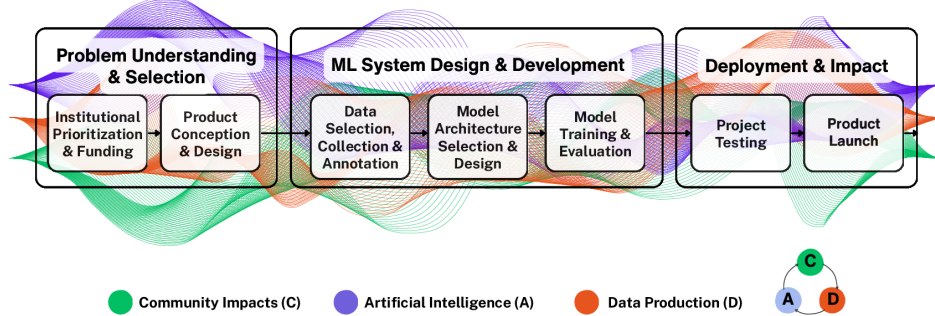


Fig. 2. Simplified AI development pipeline with three interwoven sine waves representing AI (purple), data production (orange), and community impacts (green) flowing continuously through all pipeline stages. Pipeline diagram derived from Martin [98].

those ecosystems. Parallel developments in Middle Eastern, Southeast Asian, and Latin American contexts appear less frequently, not because such initiatives were absent, but because they circulated in networks less accessible to our inquiry. We acknowledge that English-language search restrictions inherently reinforce Western-centric representation. Scholarly and community infrastructures condition what becomes visible in review corpora, creating unevenness despite our efforts. Therefore, we believe it is important to consider our own identities alongside the analysis of this work given that our backgrounds and perspectives may bias the interpretation of this work [36].

The authors hold diverse racial and ethnic identities (Black, White, Mixed-race) with cultural roots in the United States, Canada, South Africa, Ghana, Japan, and France. These backgrounds shaped our ability to recognize and access specific community-led networks (particularly in African and Indigenous contexts) while leaving others less visible to us. In terms of epistemic lens, we work across industry and academia, with backgrounds in computer science, HCI, and the humanities. This dual positioning allowed us to bridge the gap between technical documentation and critical theory, for example, recognizing “grey literature” as rigorous evidence of high-agency alternatives. However, our location within these professionalized research institutions also means we likely missed grassroots resistance tactics that do not circulate in written or digital forms. We recognize that we are observing extraction from within the institutions that often facilitate it, and we present these findings as a necessary, though partial, mapping of the landscape.

## 4 Findings

We structure our findings according to five domains of data production, which we have conceptualized as analytic elements. Rather than logistical steps, these domains function as sites where power is negotiated, contested, and encoded: **Data Relations** (the negotiation of agency and terms of engagement), **Data Labor** (the creation versus capture of value), **Data Representation** (the exercise of epistemic authority through categorization), **Data Infrastructure** (the allocation of capacity and provenance), and **Data Governance** (the enforcement of sovereignty and accountability). Within

each domain, we identify specific extractive mechanisms—technical or institutional habits that centralize control—and contrast them with high-agency pathways where communities are actively reclaiming authority.

Although mechanisms associated with each domain can appear at multiple points in AI development, consistent tendencies emerge across the corpus: decisions about relations often arise upstream as problems are framed; labor arrangements cluster within mid-stream annotation workflows; representational decisions crystallize where ontologies and preprocessing pipelines are defined; infrastructural conditions span stages but become most visible as systems scale; and governance concerns intensify downstream as models move toward evaluation and deployment. These tendencies help situate each domain without implying a fixed or linear pipeline.

Each domain is introduced through a small set of examples that surface the mechanisms we observed across the corpus. These examples are intended as points of entry into a broader landscape. The larger set of mappings, summaries, and domain categorizations is available in the datasheet for readers who wish to trace these patterns in greater depth.

#### 4.1 Data Relations

Data relations define the structural terms of engagement between model developers and the communities from whom knowledge is derived. In mainstream industry discourse, these engagements are frequently reduced to legalistic questions of copyright compliance or static “terms of service.” However, our corpus reveals that these legal frameworks often serve to obscure the underlying power dynamics [122]. Relations are not merely contractual; they are the primary site where agency is either stripped or substantiated. In extractive regimes, relations are characterized by the severance of ties between data and its creators; as Leanne Betasamosake Simpson articulates, “extraction removes all of the relationships that give whatever is being extracted meaning” [81]. High-agency relations, conversely, position data production as a negotiated partnership where community authority persists even after data is collected.

**The assumption of availability** constitutes the primary mechanism of extractive relations. Technical workflows for foundation models frequently operate on the premise that any data accessible on the public web is a “standing reserve” available for ingestion. This logic converts public existence into implicit consent. Large-scale scraping initiatives, such as the corpora used to train models like CLIP [144, 145] or T5 [131], for example, bypass the negotiation of relationship entirely, including legally, by treating the act of publication as a forfeiture of rights [77, 141]. Relation-less forms of data production systematically ignore the contextual intent of the data creator, whether it be repurposing religious texts or intimate narratives as generic linguistic tokens, none of which are “just data” [65]. By removing the requirement to ask, the assumption of availability structurally precludes the possibility of refusal, rendering the relationship unilateral.

**Transactional asymmetry** reinforces this extraction by decoupling value generation from risk. This manifests in “digital extractivism,” where Global Majority communities provide the raw material while the risks—such as the loss of privacy or the commodification of cultural heritage—are externalized back to them [69]. The dynamic functions through “accumulation by dispossession,” where the terms of engagement are dictated by the extractor, treating communities as resources rather than partners [163]. Relational asymmetries are both economic and epistemic. AI developers gain a model of the world, while communities lose control over how they are represented within it, often leading to “opportunity loss” where resources are withheld based on extractive profiling [149].

**High-agency relations** counter these mechanisms by shifting from static terms of service to dynamic and revocable consent. Rather than viewing consent as a one-time gatekeeping mechanism, high-agency approaches frame it as an ongoing relationship. The Speech Accessibility Project [2] and other initiatives that engage disability communities aptly demonstrate how relationships can precede collection: communities are partners who co-define the terms of engagement before recruiting paid volunteers and help ensure the protocol aligns with community safety needs

[2, 120]. Similarly, feminist frameworks for “embodied consent” argue for agreements that are specific, enthusiastic, and revocable, challenging the broad permissions usually buried in click-through agreements [154].

In Indigenous contexts, high-agency relations manifest as relational sovereignty. Te Hiku Media’s approach to Māori data rejects the concept of open-source availability in favor of whanaungatanga (connection/relationship), where data access is determined by the strength and trust of the relationship between the parties [29, 62]. This reintroduces friction into the data pipeline by design: access is not a default state but a negotiated privilege that requires maintaining a relationship with the originating community [83]. By replacing the assumption of availability with permissioned access, these models force a structural acknowledgment of community agency.

**Key takeaway:** Data relations determine the flow of agency. Extractive mechanisms rely on the assumption of availability, treating public data as a resource to be mined and severing the link between creators and their data. This creates transactional asymmetry, where developers capture value while communities bear the risk. High-agency relations replace this with dynamic consent and relational sovereignty, ensuring that data production remains a negotiated partnership where community authority persists throughout the technical lifecycle.

## 4.2 Data Labor

Data labor encompasses the human energy and interpretive judgment required to bridge the gap between raw information and computational capability. While often obscured by the metaphor of “autonomous” AI, our corpus confirms that model performance remains strictly dependent on human workers who select, annotate, validate, and moderate content [27, 56, 88, 107]. In extractive regimes, this labor is characterized by value capture, where the semantic value generated by human judgment is stripped from the worker and concentrated in the model, often leaving the contributor with little to no recognition or economic return. High-agency approaches, conversely, frame labor as expertise, positioning annotators as skilled contributors whose situated knowledge is essential to system quality.

**Invisibilization** by design constitutes the primary mechanism of labor extraction. Dataset and platform architectures are frequently designed to present corpora as neutral technical artifacts rather than products of human judgment, masking the interpretive decisions embedded in every labeled example [99]. This structural opacity serves to commodify the worker; by decomposing complex cultural tasks into fragmented “microtasks,” platforms strip the work of its context, rendering the worker interchangeable and the labor invisible [41]. Here, a structural design choice renders the human contribution indistinguishable from the system’s output, with the upshot of systematically preventing workers from asserting authorship claims or contesting the terms of their participation.

**Reciprocity failure** reinforces this dynamic by extracting labor without returning value. This manifests most clearly in “unwitting” labor, where user interactions (e.g., solving CAPTCHAs, tagging photos, or correcting autocomplete suggestions) are harvested to train models without the user’s explicit knowledge or compensation [18, 105]. In Global Majority contexts, this mechanism appears in the outsourcing of trauma-inducing content moderation or complex annotation to workers in low-income regions, who perform essential semantic labor for wages that do not reflect the cognitive intensity of the work [166]. The system is optimized to externalize the costs of dataset construction to the worker while centralizing the economic benefits.

**High-agency labor** counters these mechanisms by restructuring the economic and attributional relationship between modelers and workers. These approaches restore context and visibility to the labor process. The organization Karya, for example, demonstrates how data collection can function as a tool for economic redistribution; by establishing

ethical wage floors and data ownership structures for rural Indian workers, they reframe annotation as a skilled, compensated profession [3]. Similarly, the Masakhane community creates participatory research models where African language speakers function not as passive data subjects, but as credited authors and technical collaborators throughout the pipeline [112]. Emerging initiatives like Ubuntu-AI attempt to encode these rights directly into the data lifecycle through profit-sharing mechanisms, ensuring that artists and creators retain a stake in the value their data generates [111].

**Key takeaway:** Data labor is economic and political. Extractive mechanisms rely on invisibilization by design, decomposing expert judgment into fragmented tasks to obscure the worker and facilitate value capture. This severs the link between labor and downstream value. High-agency approaches replace this with labor as expertise, ensuring that contributions are visible, attributed, and compensated as skilled work that persists within the technical system.

### 4.3 Data Representation

Data representation determines how communities become computationally visible. Representation is as much a question of inclusion ratios or diversity statistics as it is one of epistemic authority: who gets to define the categories, taxonomies, and labels that structure the digital world. In extractive regimes, representation creates visibility without power, often flattening complex, relational identities into rigid categories that facilitate control or consumption. High-agency approaches, conversely, frame representation as plural epistemologies, ensuring that data structures reflect community worldviews rather than forcing local knowledge into universalizing boxes.

**Ontological imposition** constitutes the primary mechanism of representational extraction. Institutional problem formulation often imposes external taxonomies on communities before they even enter the pipeline. This manifests as “data universalism” [102] where Western logics of property and individualism are treated as neutral defaults, overwriting Indigenous ontologies that emphasize relationality and collective stewardship [89]. For example, psychological frameworks developed in “WEIRD” (Western, Educated, Industrialized, Rich, and Democratic) contexts fail to map onto collective ontologies, yet are deployed globally as standard [42, 101, 104]. Consequently, even when diverse data is collected it is structurally distorted to fit the model’s worldview, rendering specific cultural meanings “absent” even within inclusion efforts [13].

**Context stripping** reinforces this dynamic during annotation and processing. To make data “model-ready,” complex human experiences must be converted into discrete labels. This process often relies on “lazy” data practices that collapse distinct protected attributes like race and ethnicity into coarse categories to satisfy technical constraints, erasing intersectional realities [150]. Annotation workflows that lack community-defined criteria force workers to resolve ambiguity by falling back on institutional defaults, which appear neutral but encode specific cultural biases [140]. Automated filtering pipelines compound this by removing content that signals non-normative identities under the guise of “cleaning,” disproportionately purging data from non-Western contexts or disability communities [96].

**Synthetic displacement** introduces a new mechanism of extraction: representation without presence. As privacy regulations tighten, developers increasingly turn to synthetic data (e.g., fabricated medical records, artificial faces, and simulated identities) to populate datasets. While this bypasses the need for individual consent [87], it severs the link between representation and reality. Communities become represented in systems they never participated in, inheriting

the risks of misidentification or caricature without any pathway to contest how they are depicted [170]. The resulting “diversity-washing” effect is such that models appear inclusive while structurally excluding actual community members.

**High-agency representation** counters these mechanisms by building pluralistic and community-grounded corpora. These initiatives prioritize depth and context over scale. For instance, the Abundant Intelligences project reimagines AI development through Indigenous knowledge systems, refusing to separate data from the land and relations that generate it [90]. Similarly, examples from Africa and Oceania demonstrate how regional collaborations can curate datasets that serve local linguistic needs—such as the InkubaLM model—rather than adapting to global benchmarks [45, 165]. By maintaining representational authority, these projects ensure that visibility serves community goals, such as language revitalization, rather than external commodification.

**Key takeaway:** Data representation is epistemic and political. Extractive mechanisms rely on ontological imposition and context stripping, imposing external taxonomies and flattening meaning to fit technical defaults. This treats visibility as neutral even when it creates exposure. High-agency approaches replace this with plural epistemologies, grounding representation in community-defined categories and preserving the specificity of local knowledge against universalizing standards.

#### 4.4 Data Infrastructure

Data infrastructure allocates capacity and determines where data lives, who controls access, and how material circulates across model pipelines. While often treated as neutral plumbing designed for efficiency, infrastructure emerges in the literature as a primary site of political contestation. In extractive regimes, infrastructure is configured to maximize velocity and volume, creating technical conditions where consent and context are structurally impossible to maintain. High-agency approaches, conversely, design for traceability and distribution, ensuring that community authority travels with the data.

**Centralization without governance** configures extraction at an industrial scale. Foundation-model development relies on automated pipelines that ingest content from large-scale web sources such as Common Crawl or LAION to maximize throughput [43, 144, 145]. This configuration privileges actors with substantial compute resources and treats data availability as a default condition. The asymmetry is infrastructural: collection mechanisms operate at speeds that make oversight and contestation structurally unworkable for data subjects [171].

**Benchmark infrastructures** act as gatekeeping mechanisms that enforce dominant (Western) epistemologies as universal standards. Reliance on a narrow set of legacy datasets, such as ImageNet [40] and MS COCO [93], entrenches specific linguistic, cultural, and demographic assumptions as infrastructural norms [41, 82]. Because creating culturally specific alternatives requires substantial institutional support, Euro-American category systems persist as de facto standards through infrastructural path dependence [82].

**Provenance compression** serves as the third mechanism, severing datasets from their originating communities and the relational contexts of their creation. Contemporary web-scrape datasets often operate through severe documentation gaps—reinforcing “web-as-platform” assumptions that treat public accessibility as permission to extract [141]. Infrastructure that treats provenance as optional enables downstream actors to shift responsibility for data quality and rights onto untraceable contributors [94].

**High-agency infrastructure** counters these mechanisms by embedding community-defined constraints directly into technical architectures. Federated and distributed systems shift authority by enabling collaboration without

centralizing data. Emerging frameworks for “data spaces” allow communities to retain local control over storage and access while supporting model development [50, 66]. Similarly, stewardship-based architectures like Masakhane’s distributed research platforms operationalize co-designed metadata standards, ensuring that data does not become “loose” but remains tethered to its community of origin [112].

**Key takeaway:** Data infrastructure is about capacity and provenance. Extractive architectures rely on centralization without governance to maximize velocity and provenance compression to sever data from its originating obligations. This makes extraction structurally easy and accountability expensive. High-agency alternatives deploy federated and distributed systems, redistributing capacity so that community authority remains technically enforceable as data circulates.

#### 4.5 Data Governance

Governance establishes the rule-sets that authorize data production: it determines when collection is legitimate, what contextual grounding is required, and who holds authority over circulation. These rules operate upstream of participation, labor, and representation, guiding the conditions under which data production becomes legitimate. High-asymmetry governance frameworks create wide discretionary space for extractive practices, whereas high-agency governance embeds community control directly into the structures that shape data lifecycles.

**Regulatory arbitrage** constitutes the primary mechanism of extractive governance. Often termed “ethics dumping,” this practice exploits fragmented global regulations to harvest data in regions with weaker protections, converting behavioral interactions into institutional assets without oversight [161]. This dynamic transforms regulatory variation into a resource for extraction: vulnerable populations in low- and middle-income countries may receive limited digital services (like Facebook’s Free Basics) in exchange for extensive, uncompensated data harvesting [113]. Coercive collection in humanitarian settings, such as biometric registration in Ethiopian refugee camps, further illustrates how governance gaps allow institutions to bypass the consent standards required in their home jurisdictions [162].

**Open-loop extraction** reinforces this asymmetry by decoupling deployment from accountability. Models trained on narrow, Western-centric data are frequently deployed globally, shifting the burden of performance failures—such as diagnostic errors in healthcare AI—onto underserved communities [8, 115]. This mechanism externalizes risk: communities excluded from the governance of training data nonetheless become sources of performance feedback during deployment. Their interactions refine the system, yet they possess no authority to challenge the model’s adequacy or recall the data they generate [164]. Governance here functions to protect the model developer’s intellectual property while leaving the data subject’s sovereignty unprotected.

**High-agency governance** counters these mechanisms through sovereignty-based licensing and critical refusal. Rather than relying on open-access defaults, these approaches encode community authority into the legal terms of the data itself. The Kaitiakitanga License, developed by Te Hiku Media, exemplifies this by legally binding data usage to Māori tikanga (protocols), preventing extractive reuse by third parties [160]. Similarly, the Esethu Framework for African language data establishes sovereignty provisions that mandate community benefit-sharing and protect annotators [132]. Beyond licensing, critical refusal operates as a form of affirmative governance. By setting ex-ante boundaries on participation, communities assert that unreadability is a safety condition. Longstanding tactics of opacity and masking establish practical limits on what institutions may extract [22]. When viewed as governance, refusal is not a lack of data; it is an enforcement of sovereignty that limits extractive reach by design [51].

💡 **Key takeaway:** Data governance distinguishes accountability from exploitable discretion. Extractive mechanisms rely on regulatory arbitrage (ethics dumping) and open-loop extraction, engaging in collection without contextual grounding and turning deployment into unconsented data acquisition. High-agency approaches establish sovereignty-based licensing and critical refusal, creating enforceable preconditions that align data production with community-defined control and embed agency beyond the point of collection.

## 5 Discussion

Our analysis of 350 sources across academic and grey literature reveals that AI data production is not merely a logistical preliminary to model development, but a distinct sociotechnical site where power is negotiated, contested, and encoded. By synthesizing evidence across the A/D/C framework, we identify a clear divergence: extractive practices that prioritize scale, opacity, and labor externalization, versus high-agency pathways that prioritize relationality, sovereignty, and context. Notably, the five analytic domains we identify do not distribute evenly across the ML pipeline. Instead, the sources that comprise each domain cluster around the structural moments where key mechanisms take effect: **Data Relations** concentrates upstream in problem formulation and data selection; **Data Labor** anchors mid-pipeline annotation and enrichment; **Data Representation** spans early- to mid-pipeline ontology and preprocessing; **Data Infrastructure** forms a cross-cutting substrate most visible in mid-to-downstream development; and **Data Governance** clusters downstream where deployment, accountability, and sovereignty become salient. This patterned distribution indicates that extractive dynamics are not random but structurally embedded within distinct, yet interrelated, pipeline junctures.

For the HCI community, these findings suggest a critical reframing. While HCI has successfully interrogated downstream AI interaction (how users experience models) and mid-stream model behavior (bias and fairness), the upstream processes of data creation remain undertheorized in design venues. Below, we discuss how HCI scholars and practitioners can operationalize high-agency practices by treating data production as a primary site of design intervention. In doing so, we extend the nascent HCI scholarship that examines data production through the lens of data laborers and data subjects [75, 76, 100, 139, 142, 167].

### 5.1 Reframing Data Production as “Upstream” Design

Our findings challenge the industry norm of treating data as “found” infrastructure (Era 3). Instead, the evidence suggests data production is a series of design decisions—regarding relations, labor, and representation—that are often irreversible once encoded into a model. This recognition prompts us to argue that the “user” in human-centered AI must expand to include the data contributor—the artist, the annotator, the community member—whose agency is often circumvented by upstream infrastructure. Within HCI, this circumvention has predominantly been investigated in studies of data labor, which reveal how data annotators are frequently reduced to an interchangeable resource thereby constraining their subjectivities and interpretive work [75, 100, 167]. Building on this work, our review makes clear the various design choices like crowdsourcing interfaces that atomize tasks to obscure the worker’s context (Data Labor) or scraping pipelines that strip provenance metadata (Data Infrastructure) which can circumvent agency and enforce extraction by design.

For HCI, this implies that data curation is a form of interaction design. The high-agency pathways our review surfaces make clear that alternative designs are possible. The community-led initiatives such as Masakhane’s participatory NLP

[97] or Māori data sovereignty protocols [83, 108] succeed not by “fixing” extraction after the fact, but by designing relational friction into the process. They replace the seamless, frictionless extraction of web scraping with protocols that require consent, negotiation, and maintenance.

## 5.2 Toward High-Agency Practices: Implications for HCI

Moving beyond critique, our analysis of high-agency pathways points toward concrete mechanisms for less-extractive AI development. We map these implications to three key shifts for HCI practice.

*5.2.1 From Universal Representation to Pluralistic “Small Data”.* The dominance of massive, web-scraped corpora (Era 3) enforces a universalizing worldview that erases minority contexts and agency in general. Our findings suggest that “de-biasing” these massive datasets is often less effective than building smaller, community-sovereign corpora. Indeed, critiques of contemporary efforts to build more “inclusive” or “de-biased” technologies often highlight a technosolutionist trap whereby the issue is purportedly addressed through large-scale capture of data about a community or culture without meaningful agency [26, 129, 130]. Our review reinforces how failing to allow communities to set the terms of inclusion for their data can inadvertently perpetuate extraction under the guise of inclusion. This extends critiques of “fair” AI that don’t fundamentally shift power [73]. In contrast, high-agency pathways highlighted in our review demonstrate how different actors and groups are proactively responding to extractive data capture by imagining and building alternatives. For instance, efforts by Te Hiku Media to develop community led language datasets and technologies [45, 160] and the Community-driven African Next Voices project [176] directly counter efforts by big tech to seek out and capture cultural knowledge and data, by instead keeping data governed by communities. While authors, organizations, and initiatives offer unique contributions, their coordination and totality points to systemic alternatives that start with community needs and maintain community control, thus creating their own conditions for thriving rather than adapting to external constraints. By developing their own evaluation criteria, publication venues, funding mechanisms, and governance protocols, they establish parallel infrastructures that operate according to different principles: sovereignty rather than extraction, reciprocity rather than accumulation, cultural preservation rather than homogenization and standardization.

The HCI community has an opportunity to advance high-agency efforts by investing in federated data spaces rather than centralized lakes. We need infrastructure that allows models to learn from community data without that data ever leaving the community’s local storage or jurisdiction (e.g., federated learning tailored for Indigenous sovereignty [50]). Furthermore, valuation metrics in AI research must shift away from scale at the expense of care [61, 71, 139, 175]. We encourage the HCI community to value (and publish) contributions that curate high-context, small-scale datasets with clear governance protocols, rather than rejecting them for lacking the scale of foundation model benchmarks.

*5.2.2 From Transactional Labor to Relational Provenance.* Our review highlights reciprocity failure and labor invisibility as central extractive patterns. This transactional model commodifies the work of data laborers—including annotators, content creators, and community members—distancing those performing the work from those capturing the value and contributing to the perceived “magic” of AI [137]. This is further complicated by opaque data collection practices that often make the data creator unaware of their contribution. While there is growing interest in data provenance as a key intervention point for mitigating harm of AI technologies [94, 169], our review affirms how the current dominant paradigm of data production is in tension with this end goal. We argue that this tension is, in part, a design challenge for HCI researchers and call upon the community to explore how data capture and sharing platforms might implement provenance-tracking mechanisms by design.

A provenance-first design approach could involve binding labor and authorship metadata to individual data points so that creators retain “credits” (similar to the Ubuntu-AI model [111]) that persist through the pipeline. More broadly, there is growing recognition that data annotation is fundamentally subjective and interpretive, often shaped by the sociocultural backgrounds and lived experiences of annotators [46, 146? ]. This motivates the design of data annotation processes and infrastructure that allows workers to signal ambiguity, refuse tasks that violate community norms, and capture disagreements in a structured form rather than forcing a choice that flattens cultural context. By doing so, the HCI community can enable downstream pluralistic modeling approaches that can handle meaningful divergences in perspectives [106].

**5.2.3 From Open-Loop Extraction to Closed-Loop Governance.** The governance gaps identified in our review show that once data is scraped, communities often lose control. In contrast, community-led governance approaches exemplify an alternative. For example, Māori data sovereignty frameworks in Aotearoa New Zealand demonstrate a coordinated ecosystem: the Māori Data Sovereignty Network develops governance protocols, Te Hiku Media creates community-led, culturally appropriate datasets and benchmarks, and the Kaitiakitanga License embeds community authority into legal frameworks. Such agency-oriented practices require dynamic consent and enforceable boundaries and necessitate technical implementations of Sovereignty-Based Licensing. HCI scholars can develop and standardize machine-readable licenses (similar to Creative Commons but for ML training) that explicitly forbid certain downstream uses (e.g., military application, generative mimicry) and trigger benefit-sharing clauses [132, 160].

To operationalize this, Institutional Review Boards (IRBs) and conference ethics reviews must look upstream, enforcing data transparency standards that treat data collection as a distinct object of ethical inquiry [11, 54]. Within academic peer review, data production is increasingly within scope of ethical inquiry. However, the focus remains largely on individual privacy and consent within papers presenting novel datasets, rather than deeper inquiries into the conditions under which data is produced and the extent to which communities retain any rights of refusal. This necessitates new review paradigms that prioritizes community consent in addition to individual terms of service, echoing calls for power-aware approaches that allows communities to attest and refuse data extraction [51, 73]. By scrutinizing these data cascades at the source [136], the review process can identify where the agency of the data contributor has been circumvented.

### 5.3 Methodological Contributions: The Value of Multivocality

Finally, this paper validates the utility of the Multivocal Literature Review (MLR) for investigating sociotechnical harm. A significant portion of our high-agency evidence came not from peer-reviewed academic venues, but from “grey” literature—community manifestos, tribal resolutions, and worker inquiries. If we had limited our scope to academic “white” literature, we would have successfully diagnosed the harms of extraction (which are well-documented in academia) but missed the existing alternatives (which are often documented in policy and community organizing). For the HCI community, this underscores that “state-of-the-art” knowledge regarding justice and equity often resides outside the academy, as does the general “state-of-practice.” Future work on AI harms should adopt multivocal methods to ensure that community-generated resistance and innovation are recognized as rigorous evidence.

## 6 Conclusions

As AI development consolidates around foundation models trained on internet-scale scrapes, the risk of deepening extractive relations is acute. However, this trajectory is not inevitable. By analyzing the data production pipeline

through the lens of Data Relations, Data Labor, Data Representation, Data Infrastructure, and Data Governance, we see that every dataset is a record of power relations. This paper contributes a taxonomy of these relations, offering HCI a diagnostic tool to identify extraction and a catalog of precedents for resistance. The shift to less-extractive AI requires more than better algorithms; it requires designing the upstream sociotechnical infrastructures that determine whose knowledge counts, how it is valued, and who governs its future. Our review affirms that a less-extractive future is not merely an aspiration; it is actively being built by communities pursuing alternative to the status quo.

## References

- [1] [n. d.]. Data Workers Inquiry. <https://data-workers.org/>
- [2] [n. d.]. Speech Accessibility Project. <https://speechaccessibilityproject.beckman.illinois.edu>
- [3] Basil Abraham, Danish Goel, Divya Siddarth, Kalika Bali, Manu Chopra, Monojit Choudhury, Pratik Joshi, Preethi Jyoti, Sunayana Sitaram, and Vivek Seshadri. 2020. Crowdsourcing Speech Data for Low-Resource Languages from Low-Income Workers. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (Eds.). European Language Resources Association, Marseille, France, 2819–2826. <https://aclanthology.org/2020.lrec-1.343/>
- [4] Mark Adley, Hayley Alderson, Katherine Jackson, William McGovern, Liam Spencer, Michelle Addison, and Amy O'Donnell. 2024. Ethical and practical considerations for including marginalised groups in quantitative survey research. *International Journal of Social Research Methodology* 27, 5 (2024), 559–574. doi:10.1080/13645579.2023.2228600
- [5] Leah Hope Ajmani, Jasmine C. Foriest, Jordan Taylor, Kyle Pittman, Sarah Gilbert, and Michael Ann DeVito. 2024. Whose Knowledge is Valued? Epistemic Injustice in CSCW Applications. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW2 (Nov. 2024), 523:1–523:28. doi:10.1145/3687062
- [6] Hilary Arksey and Lisa O'Malley. 2005. Scoping studies: towards a methodological framework. *International Journal of Social Research Methodology* 8, 1 (Feb. 2005), 19–32. doi:10.1080/1364557032000119616
- [7] A. Arora, M. Barrett, E. Lee, E. Oborn, and K. Prince. 2023. Risk and the future of AI: Algorithmic bias, data colonialism, and marginalization. *Information and Organization* 33, 3 (2023), 100478. doi:10.1016/j.infoandorg.2023.100478
- [8] Mercy Asiedu, Awa Dieng, Iskandar Haykel, Negar Rostamzadeh, Stephen Pfohl, Chirag Nagpal, Maria Nagawa, Abigail Oppong, Sanmi Koyejo, and Katherine Heller. 2024. The Case for Globalizing Fairness: A Mixed Methods Study on Colonialism, AI, and Health in Africa. arXiv:2403.03357 (March 2024). doi:10.48550/arXiv.2403.03357 arXiv:2403.03357 [cs].
- [9] Stefan Baack. 2024. A Critical Analysis of the Largest Source for Generative AI Training Data: Common Crawl. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 2199–2208. doi:10.1145/3630106.3659033
- [10] Marie Battiste. 2005. Indigenous Knowledge: Foundations for First Nations. *Worm Indigenous Nations Higher Education Consortium Journal* (Jan. 2005). [https://www.researchgate.net/publication/241822370\\_Indigenous\\_Knowledge\\_Foundations\\_for\\_First\\_Nations](https://www.researchgate.net/publication/241822370_Indigenous_Knowledge_Foundations_for_First_Nations)
- [11] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics* 6 (Dec. 2018), 587–604. doi:10.1162/tac1\_a\_00041
- [12] Eshta Bhardwaj, Harshit Gujral, Siyi Wu, Ciara Zogheib, Tegan Maharaj, and Christoph Becker. 2024. Machine learning data practices through a data curation lens: An evaluation framework. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 1055–1067. doi:10.1145/3630106.3658955
- [13] Steven Bird. 2024. Must NLP be Extractive? [https://drive.google.com/file/d/1hvF7\\_WQrou6CWZydhymYFTYHnd3ZljV/view?usp=embed\\_facebook](https://drive.google.com/file/d/1hvF7_WQrou6CWZydhymYFTYHnd3ZljV/view?usp=embed_facebook)
- [14] Abeba Birhane. 2020. Algorithmic Colonization of Africa. *SCRIPTed* 17, 2 (Aug. 2020), 389–409. doi:10.2966/scrip.170220.389
- [15] Abeba Birhane. 2021. Algorithmic injustice: a relational ethics approach. *Patterns* 2, 2 (Feb. 2021), 100205. doi:10.1016/j.patter.2021.100205
- [16] Abeba Birhane, William Isaac, Vinodkumar Prabhakaran, Mark Diaz, Madeleine Clare Elish, Iason Gabriel, and Shakir Mohamed. 2022. Power to the People? Opportunities and Challenges for Participatory AI. In *Equity and Access in Algorithms, Mechanisms, and Optimization*. ACM, Arlington VA USA, 1–8. doi:10.1145/3551624.3555290
- [17] Abeba Birhane, Elaine Ruane, Thomas Laurent, Matthew S. Brown, Johnathan Flowers, Anthony Ventresque, and Christopher L. Dancy. 2022. The Forgotten Margins of AI Ethics. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*. Association for Computing Machinery, New York, NY, USA, 948–958. doi:10.1145/3531146.3533157
- [18] Briony Blackmore, Michelle Thorp, Andrew Tzer-Yeu Chen, Fabio Morreale, Brent Burmester, Elham Bahmanteymouri, and Matt Bartlett. 2023. Hidden humans: exploring perceptions of user-work and training artificial intelligence in Aotearoa New Zealand. *Kōtuitui: New Zealand Journal of Social Sciences Online* 18, 4 (Oct. 2023), 443–456. doi:10.1080/1177083X.2023.2212736
- [19] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of “Bias” in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 5454–5476. doi:10.18653/v1/2020.acl-main.485
- [20] John R. Bowman. 1989. *Capitalist Collective Action: Competition, Cooperation and Conflict in the Coal Industry*. Cambridge University Press. Google-Books-ID: fnl6sAYpRLYC.
- [21] Paul T. Brown, Daniel Wilson, Kiri West, Kirita-Rose Escott, Kiya Basabas, Ben Ritchie, Danielle Lucas, Ivy Taia, Natalie Kusabs, and Te Taka Keegan. 2024. Māori Algorithmic Sovereignty: Idea, Principles, and Use. *Data Science Journal* 23, 1 (April 2024). doi:10.5334/dsj-2024-015
- [22] Simone Browne. 2015. *Dark Matters: On the Surveillance of Blackness*. Duke University Press. doi:10.1215/9780822375302
- [23] Jenna Burrell. 2016. How the machine ‘thinks’: Understanding opacity in machine learning algorithms. *Big Data & Society* 3, 1 (2016), 2053951715622512. doi:10.1177/2053951715622512
- [24] Judith Butler. 1990. *Gender Trouble: Feminism and the Subversion of Identity*. Routledge. Google-Books-ID: kuztAAAAAAJ.
- [25] Stephanie Russo Carroll, Pyrou Chung, Robyn K. Rowe, Susanna Siri, and Walter. 2024. Indigenous Data Sovereignty and the State of Open Data. <https://www.d4d.net/news/indigenous-data-sovereignty-and-the-state-of-open-data>

- [26] Alan Chan, Chinsaa T Okolo, Zachary Turner, and Angelina Wang. 2021. The limits of global inclusion in AI development. *arXiv preprint arXiv:2102.01265* (2021).
- [27] Srravya Chandhiramowuli, Alex S. Taylor, Sara Heitlinger, and Ding Wang. 2024. Making Data Work Count. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW1 (April 2024), 90:1–90:26. doi:10.1145/3637367
- [28] Ishita Chordia, Leya Breanna Baltaxe-Admony, Ashley Boone, Alyssa Sheehan, Lynn Dombrowski, Christopher A Le Dantec, Kathryn E. Ringland, and Angela D. R. Smith. 2024. Social Justice in HCI: A Systematic Literature Review. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–33. doi:10.1145/3613904.3642704
- [29] Donavyn Coffey. 2021. Māori are trying to save their language from Big Tech. *Wired* (April 2021). <https://www.wired.com/story/maori-language-tech/>
- [30] Cathy J. Cohen. 1997. Punks, Bulldaggers, and Welfare Queens: The Radical Potential of Queer Politics? *GLQ: A Journal of Lesbian and Gay Studies* 3, 4 (1997), 437–465.
- [31] Combahee River Collective. 1977. (1977) The Combahee River Collective Statement •. <https://www.blackpast.org/african-american-history/combahee-river-collective-statement-1977/>
- [32] Patricia Hill Collins. 2002. *Black Feminist Thought: Knowledge, Consciousness, and the Politics of Empowerment* (2 ed.). Routledge, New York. doi:10.4324/9780203900055
- [33] Common Crawl. 2025. *Common Crawl*. <https://commoncrawl.org/>
- [34] Ned Cooper, Tiffanie Horne, Gillian R Hayes, Courtney Heldreth, Michal Lahav, Jess Holbrook, and Lauren Wilcox. 2022. A Systematic Review and Thematic Analysis of Community-Collaborative Approaches to Computing Research. In *CHI Conference on Human Factors in Computing Systems*. ACM, New Orleans LA USA, 1–18. doi:10.1145/3491102.3517716
- [35] Matthew Cotton. 2017. Fair fracking? Ethics and environmental justice in United Kingdom shale gas policy and planning. *Local Environment* 22, 2 (Feb. 2017), 185–202. doi:10.1080/13549839.2016.1186613
- [36] Payton Croskey, Fabian Offert, Jennifer Jacobs, and Kai M. Thaler. 2025. Liberatory Collections and Ethical AI: Reimagining AI Development from Black Community Archives and Datasets. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency (FAccT '25)*. Association for Computing Machinery, New York, NY, USA, 900–913. doi:10.1145/3715275.3732058
- [37] Íñigo de Troya, Jacqueline Kernahan, Neelke Doorn, Virginia Dignum, and Roel Dobbe. 2025. Misabstraction in Sociotechnical Systems. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency (FAccT '25)*. Association for Computing Machinery, New York, NY, USA, 1829–1842. doi:10.1145/3715275.3732122
- [38] Jeffrey Dean. 2019. The Deep Learning Revolution and Its Implications for Computer Architecture and Chip Design. arXiv:1911.05289 (Nov. 2019). doi:10.48550/arXiv.1911.05289 arXiv:1911.05289 [cs].
- [39] Vine Deloria and Clifford M. Lytle. 1998. *The Nations Within: The Past and Future of American Indian Sovereignty*. University of Texas Press. Google-Books-ID: FLgEf5kGLWQC.
- [40] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, Miami, FL, 248–255. doi:10.1109/CVPR.2009.5206848
- [41] Emily Denton, Alex Hanna, Razvan Amironesei, Andrew Smart, and Hilary Nicole. 2021. On the genealogy of machine learning datasets: A critical history of ImageNet. *Big Data & Society* 8, 2 (2021), 205395172110359. doi:10.1177/20539517211035955
- [42] Lindsey DeWitt Prat, Olivia Nercy Ndlovu Lucas, Christopher Golias, and Mia Lewis. 2024. Decolonizing LLMs: An Ethnographic Framework for AI in African Contexts. *EPIC Proceedings* (2024), 45–84. <https://doi.org/10.1111/epic.12196>
- [43] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. arXiv:2104.08758 (2021). doi:10.48550/arXiv.2104.08758 arXiv:2104.08758 [cs].
- [44] Pedro Domingos. 2012. A few useful things to know about machine learning. *Commun. ACM* 55, 10 (Oct. 2012), 78–87. doi:10.1145/2347736.2347755
- [45] Suzanne Duncan, Gianna Leoni, Lee Steven, Keoni Mahelona, and Peter-Lucas Jones. 2024. Fit for our purpose, not yours: Benchmark for a low-resource, Indigenous language. <https://openreview.net/forum?id=w5jfyvsRq3#discussion>
- [46] Mark Díaz, Ian Kivlichan, Rachel Rosen, Dylan Baker, Razvan Amironesei, Vinodkumar Prabhakaran, and Remi Denton. 2022. CrowdWorkSheets: Accounting for Individual and Collective Identities Underlying Crowdsourced Dataset Annotation. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*. Association for Computing Machinery, New York, NY, USA, 2342–2351. doi:10.1145/3531146.3534647
- [47] Pedro Ferreira. 2024. Examining the “Local” in ICT4D: A Postcolonial Perspective on Participation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–13. doi:10.1145/3613904.3642748
- [48] Raymond Fok, Alexa Siu, and Daniel S. Weld. 2025. Toward Living Narrative Reviews: An Empirical Study of the Processes and Challenges in Updating Survey Articles in Computing Research. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, 1–10. doi:10.1145/3706598.3714047
- [49] Miranda Fricker. 2007. *Epistemic injustice: power and the ethics of knowing*. Oxford university press, Oxford.
- [50] Ana García, Savvas Rogotis, Eimear Farrell, Tobias Guggenberger, Arash Hajikhani, Atte Kinnula, Marko Komssi, and Tuomo Tuikka. 2024. Generative AI and Data Spaces: White Paper. (2024).

- [51] Patricia Garcia, Tonia Sutherland, Niloufar Salehi, Marika Cifor, and Anubha Singh. 2022. No! Re-imagining Data Practices Through the Lens of Critical Refusal. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW2 (Nov. 2022), 315:1–315:20. doi:10.1145/3557997
- [52] Vahid Garousi, Michael Felderer, and Mika V. Mäntylä. 2016. The need for multivocal literature reviews in software engineering: complementing systematic literature reviews with grey literature. In *Proceedings of the 20th International Conference on Evaluation and Assessment in Software Engineering (EASE '16)*. Association for Computing Machinery, New York, NY, USA, 1–6. doi:10.1145/2915970.2916008
- [53] Vahid Garousi, Michael Felderer, and Mika V. Mäntylä. 2019. Guidelines for including grey literature and conducting multivocal literature reviews in software engineering. *Information and Software Technology* 106 (Feb. 2019), 101–121. doi:10.1016/j.infsof.2018.09.006
- [54] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets. *Commun. ACM* 64, 12 (Nov. 2021), 86–92. doi:10.1145/3458723
- [55] Katelyn Godin, Jackie Stapleton, Sharon I. Kirkpatrick, Rhona M. Hanning, and Scott T. Leatherdale. 2015. Applying systematic review search methods to the grey literature: a case study examining guidelines for school-based breakfast programs in Canada. *Systematic Reviews* 4 (Oct. 2015), 138. doi:10.1186/s13643-015-0125-0
- [56] Mary L Gray and Siddharth Suri. 2019. *Ghost Work: How to Stop Silicon Valley from Building a New Global Underclass*. Houghton Mifflin Harcourt.
- [57] Michael Haenlein and Andreas Kaplan. 2019. A Brief History of Artificial Intelligence: On the Past, Present, and Future of Artificial Intelligence. *California Management Review* 61, 4 (Aug. 2019), 5–14. doi:10.1177/0008125619864925
- [58] Alon Halevy, Peter Norvig, and Fernando Pereira. 2009. The Unreasonable Effectiveness of Data. *IEEE Intelligent Systems* 24, 2 (March 2009), 8–12. doi:10.1109/MIS.2009.36
- [59] Siobhan Mackenzie Hall, Samantha Dalal, Raesetje Sefala, Foutse Yueghog, Aisha Alaagib, Imane Hamzaoui, Shu Ishida, Jabez Magomere, Lauren Crais, Aya Salama, and Tejumade Afonja. 2025. The Human Labour of Data Work: Capturing Cultural Diversity through World Wide Dishes. arXiv:2502.05961 (Feb. 2025). doi:10.48550/arXiv.2502.05961 arXiv:2502.05961 [cs].
- [60] Jamie Hancock, Sarada Mahesh, Jennifer Cobbe, Jatinder Singh, and Anjali Mazumder. 2024. The tensions of data sharing for human rights: A modern slavery case study. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 974–987. doi:10.1145/3630106.3658949
- [61] Alex Hanna and Tina M. Park. 2020. Against Scale: Provocations and Resistances to Scale Thinking. arXiv:2010.08850 (Nov. 2020). doi:10.48550/arXiv.2010.08850 arXiv:2010.08850 [cs].
- [62] Karen Hao. 2022. A new vision of artificial intelligence for the people. *MIT Technology Review* (April 2022). <https://www.technologyreview.com/2022/04/22/1050394/artificial-intelligence-for-the-people/>
- [63] Karen Hao and Andrea Paola Hernández. 2022. How the AI industry profits from catastrophe. *MIT Technology Review* (April 2022). <https://www.technologyreview.com/2022/04/20/1050392/ai-industry-appen-scale-data-labels/>
- [64] Donna Haraway. 1988. Situated Knowledges: The Science Question in Feminism and the Privilege of Partial Perspective. *Feminist Studies* 14, 3 (1988), 575–599. doi:10.2307/3178066
- [65] Ben Hutchinson. 2024. Modeling the Sacred: Considerations when Using Religious Texts in Natural Language Processing. arXiv:2404.14740 (2024). doi:10.48550/arXiv.2404.14740 arXiv:2404.14740 [cs].
- [66] Andreas Hutterer and Barbara Krumay. 2024. The adoption of data spaces: Drivers toward federated data sharing. doi:10.24251/HICSS.2024.542
- [67] Paul Agu Igwe, Nnamdi O. Madichie, and David Gamariel Rugara. 2022. Decolonising research approaches towards non-extractive research. *Qualitative Market Research: An International Journal* 25, 4 (Jan. 2022), 453–468. doi:10.1108/QMR-11-2021-0135
- [68] Lilly Irani, Janet Vertesi, Paul Dourish, Kavita Philip, and Rebecca E. Grinter. 2010. Postcolonial computing: a lens on design and development. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Atlanta Georgia USA, 1311–1320. doi:10.1145/1753326.1753522
- [69] Neema Iyer, Garnett Achieng, Favour Borokini, Uri Ludger, Neema Iyer, Yahya Syabani, and Yahya Syabani. 2021. Automated Imperialism, Expansionist Dreams: Exploring Digital Extractivism in Africa. (2021). <https://archive.policyp.org/digitalextractivism/>
- [70] C. L. R. James. 1989. *The black Jacobins: Toussaint l'Ouverture and the San Domingo revolution* (2. ed., rev ed.). Vintage Books, a Division of Random House, Inc, New York.
- [71] Eun Seo Jo and Timnit Gebru. 2020. Lessons from archives: strategies for collecting sociocultural data in machine learning. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT\* '20)*. Association for Computing Machinery, New York, NY, USA, 306–316. doi:10.1145/3351095.3372829
- [72] James H Jones. 2008. The Tuskegee syphilis experiment. *The Oxford textbook of clinical research ethics* (2008), 86–96.
- [73] Pratyusha Kalluri. 2020. Don't ask if artificial intelligence is good or fair, ask how it shifts power. *Nature* 583, 7815 (2020), 169–169.
- [74] Fernando Kamei, Igor Wiese, Gustavo Pinto, Waldemar Ferreira, Márcio Ribeiro, Renata Souza, and Sérgio Soares. 2022. Assessing the Credibility of Grey Literature: A Study with Brazilian Software Engineering Researchers. *Journal of Software Engineering Research and Development* 10 (june 2022). doi:10.5753/jsrerd.2022.1897
- [75] Shivani Kapania, Alex S Taylor, and Ding Wang. 2023. A hunt for the Snark: Annotator Diversity in Data Practices. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. ACM, Hamburg Germany, 1–15. doi:10.1145/3544548.3580645
- [76] Reishiro Kawakami and Sukrit Venkatagiri. 2024. The Impact of Generative AI on Artists. In *Proceedings of the 16th Conference on Creativity & Cognition (Camp/C '24)*. Association for Computing Machinery, New York, NY, USA, 79–82. doi:10.1145/3635636.3664263
- [77] Mehtab Khan and Alex Hanna. 2022. The Subjects and Stages of AI Dataset Development: A Framework for Dataset Accountability. *Forthcoming 19 Ohio St. Tech. L.J. (2023)* (2022). doi:10.2139/ssrn.4217148

- [78] Rob Kitchin, Juliette Davret, Carla M Kayanan, and Samuel Mutter. 2025. Assemblage theory, data systems and data ecosystems: The data assemblages of the Irish planning system. *Big Data & Society* 12, 3 (2025), 20539517251352822. doi:10.1177/20539517251352822
- [79] Rob Kitchin and Tracey Lauriault. 2014. Towards Critical Data Studies: Charting and Unpacking Data Assemblages and Their Work. (2014). <https://papers.ssrn.com/abstract=2474112>
- [80] Lauren Klein and Catherine D'Ignazio. 2024. Data Feminism for AI. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 100–112. doi:10.1145/3630106.3658543
- [81] Naomi Klein. 2013. Naomi Klein Chats with Leanne Simpson about Idle No More. <https://www.yesmagazine.org/social-justice/2013/03/06/dancing-the-world-into-being-a-conversation-with-idle-no-more-leanne-simpson>
- [82] Bernard Koch, Emily Denton, Alex Hanna, and Jacob G. Foster. 2021. Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research. arXiv:2112.01716 (Dec. 2021). doi:10.48550/arXiv.2112.01716 arXiv:2112.01716 [cs, stat].
- [83] Tahu Kukutai and Donna Cormack. 2020. *"Pushing the space": Data sovereignty and self-determination in Aotearoa NZ* (1st edition ed.). Routledge, 21–35. <https://www.taylorfrancis.com/reader/read-online/6abf9fc2-820b-4950-b310-eef574873fbb/chapter/pdf?context=ubx>
- [84] Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Minh Chien Vu, Itziar González-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepereq, Suzana Ilic, Margaret Mitchell, Sasha Alexandra Luccioni, and Yacine Jernite. 2023. The BigScience ROOTS Corpus: A 1.6TB Composite Multilingual Dataset. arXiv:2303.03915 (March 2023). doi:10.48550/arXiv.2303.03915 arXiv:2303.03915 [cs].
- [85] Y. LeCun. [n. d.]. THE MNIST DATABASE of handwritten digits. <http://yann.lecun.com/exdb/mnist/> ([n. d.]). <https://cir.nii.ac.jp/crid/1571417126193283840>
- [86] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (Nov. 1998), 2278–2324. doi:10.1109/5.726791
- [87] Peter Lee. 2024. Synthetic Data and the Future of AI. 4722162 (Feb. 2024). <https://papers.ssrn.com/abstract=4722162>
- [88] Tuukka Lehtiniemi and Minna Ruckenstein. 2022. *Prisoners training AI: Ghosts, Humans and Values in Data Labour*. Routledge, Abingdon, 184–196. doi:10.4324/9781003170884-16
- [89] Jason Edward Lewis, Angie Abdilla, Noelani Arista, Kaipulaumakaniolono Baker, Scott Benesiinaabandan, Michelle Brown, Melanie Cheung, Meredith Coleman, Ashley Cordes, Joel Davison, Kūpono Duncan, Sergio Garzon, D. Fox Harrell, Peter-Lucas Jones, Kekuhi Kealiikanakaoleo-haililani, Megan Kelleher, Suzanne Kite, Olin Lagon, Jason Leigh, Maroussia Levesque, Keoni Mahelona, Caleb Moses, Isaac ('Ika'aka) Nahuwai, Kari Noe, Danielle Olson, 'Ōiwi Parker Jones, Caroline Running Wolf, Michael Running Wolf, Marlee Silva, Skawennati Fragnito, and Hēmi Whaanga. 2020. *Indigenous Protocol and Artificial Intelligence Position Paper*. Indigenous Protocol and Artificial Intelligence Working Group and the Canadian Institute for Advanced Research, Honolulu, HI. doi:10.11573/spectrum.library.concordia.ca.00986506
- [90] Jason Edward Lewis, Hēmi Whaanga, and Ceyda Yolgörmez. 2025. Abundant intelligences: placing AI within Indigenous knowledge frameworks. *AI & SOCIETY* 40, 4 (April 2025), 2141–2157. doi:10.1007/s00146-024-02099-4
- [91] Calvin A. Liang, Sean A. Munson, and Julie A. Kientz. 2021. Embracing Four Tensions in Human-Computer Interaction Research with Marginalized People. *ACM Trans. Comput.-Hum. Interact.* 28, 2 (April 2021), 14:1–14:47. doi:10.1145/3443686
- [92] Andreas Liesenfeld and Mark Dingemanse. 2024. Rethinking open source generative AI: open washing and the EU AI Act. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 1774–1787. doi:10.1145/3630106.3659005
- [93] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [94] Shayne Longpre, Nikhil Singh, Manuel Cherep, Kushagra Tiwary, Joanna Materzynska, William Brannon, Robert Mahari, Manan Dey, Mohammed Hamdy, Nayan Saxena, Ahmad Mustafa Anis, Emad A. Alghamdi, Vu Minh Chien, Naana Obeng-Marnu, Da Yin, Kun Qian, Yizhi Li, Minnie Liang, An Dinh, Shrestha Mohanty, Deividas Mataciunas, Tobin South, Jianguo Zhang, Ariel N. Lee, Campbell S. Lund, Christopher Klam, Damien Sileo, Diganta Misra, Enrico Shippole, Kevin Klyman, Lester JV Miranda, Niklas Muennighoff, Seonghyeon Ye, Seungone Kim, Vipul Gupta, Vivek Sharma, Xuhui Zhou, Caiming Xiong, Luis Villa, Stella Biderman, Alex Pentland, Sara Hooker, and Jad Kabbara. 2024. Bridging the Data Provenance Gap Across Text, Speech and Video. arXiv:2412.17847 (Dec. 2024). doi:10.48550/arXiv.2412.17847 arXiv:2412.17847 [cs].
- [95] Qinghua Lu, Liming Zhu, Xiwei Xu, Jon Whittle, Didar Zowghi, and Aurelie Jacquet. 2024. Responsible AI Pattern Catalogue: A Collection of Best Practices for AI Governance and Engineering. *Comput. Surveys* 56, 7 (2024), 1–35. doi:10.1145/3626234
- [96] Kelly Avery Mack, Rida Qadri, Remi Denton, Shaun K. Kane, and Cynthia L. Bennett. 2024. "They only care to show us the wheelchair": disability representation in text-to-image AI models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–23. doi:10.1145/3613904.3642166
- [97] Vukosi Marivate. 2021. *Why African natural language processing now? A view from South Africa AfricaNLP*. Mapungubwe Institute for Strategic Reflection (MISTRA), 126–152. doi:10.2307/jj.12406168.11
- [98] Donald Jr. Martin. 2020. Upgrading the Product Development Process to Foster Machine Learning Fairness and Ethical AI. <https://www.youtube.com/watch?v=1U9cSPeYkA>

- [99] Milagros Miceli and Julian Posada. 2022. The Data-Production Dispositif. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov. 2022), 1–37. doi:10.1145/3555561
- [100] Milagros Miceli, Martin Schuessler, and Tianling Yang. 2020. Between Subjectivity and Imposition: Power Dynamics in Data Annotation for Computer Vision. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (Oct. 2020), 1–25. doi:10.1145/3415186
- [101] Rada Mihalcea, Oana Ignat, Longju Bai, Angana Borah, Luis Chiruzzo, Zhijing Jin, Claude Kwizera, Joan Nwatu, Soujanya Poria, and Tamar Solorio. 2025. Why AI is WEIRD and shouldn't be this way: towards AI for everyone, with everyone, by everyone. In *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence (AAAI'25/IAAI'25/EAAI'25, Vol. 39)*. AAAI Press, 28657–28670. doi:10.1609/aaai.v39i27.35092
- [102] Stefania Milan and Emiliano Treré. 2019. Big Data from the South(s): Beyond Data Universalism. *Television New Media* 20, 4 (May 2019), 319–335. doi:10.1177/1527476419837739
- [103] Shakir Mohamed, Marie-Therese Png, and William Isaac. 2020. Decolonial AI: Decolonial Theory as Sociotechnical Foresight in Artificial Intelligence. *Philosophy Technology* 33, 4 (Dec. 2020), 659–684. doi:10.1007/s13347-020-00405-8
- [104] Cristina Jayme Montiel and Joshua Uyheng. 2022. Foundations for a decolonial big data psychology. *Journal of Social Issues* 78, 2 (2022), 278–297. doi:10.1111/josi.12439
- [105] Fabio Morreale, Elham Bahmanteymouri, Brent Burmester, Andrew Chen, and Michelle Thorp. 2023. The unwitting labourer: extracting humanness in AI training. *AI SOCIETY* (May 2023). doi:10.1007/s00146-023-01692-3
- [106] Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022. Dealing with Disagreements: Looking Beyond the Majority Vote in Subjective Annotations. *Transactions of the Association for Computational Linguistics* 10 (2022), 92–110. doi:10.1162/tacl\_a\_00449
- [107] James Muldoon, Callum Cant, Boxi Wu, and Mark Graham. 2024. A Typology of AI Data Work. *Big Data and Society* 11, 11 (March 2024). doi:10.1177/20539517241232632
- [108] Luke Munn. 2024. The five tests: designing and evaluating AI according to indigenous Māori principles. *AI SOCIETY* 39, 4 (Aug. 2024), 1673–1681. doi:10.1007/s00146-023-01636-x
- [109] Zachary Munn, Micah D. J. Peters, Cindy Stern, Catalin Tufanaru, Alexa McArthur, and Edoardo Aromataris. 2018. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping review approach. *BMC Medical Research Methodology* 18, 1 (Nov. 2018), 143. doi:10.1186/s12874-018-0611-x
- [110] Diego I. Murguía and Kathrin Böhlting. 2013. Sustainability reporting on large-scale mining conflicts: the case of Bajo de la Alumbrera, Argentina. *Journal of Cleaner Production* 41 (Feb. 2013), 202–209. doi:10.1016/j.jclepro.2012.10.012
- [111] M. Nayebare, R. Eglash, U. Kimanuku, R. Baguma, J. Mounsey, and C. Maina. 2023. *Interim Report for Ubuntu-AI: A Bottom-up Approach to More Democratic and Equitable Training and Outcomes for Machine Learning*. San Francisco. <https://generativejustice.org/uai/>
- [112] Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohunge, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Seling, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamil Toure Ali, Jade Abbott, Iroko Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elshahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. 2020. Participatory Research for Low-resourced Machine Translation: A Case Study in African Languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 2144–2160. doi:10.18653/v1/2020.findings-emnlp.195
- [113] Toussaint Nothias. 2020. Access granted: Facebook's free basics in Africa. *Media, Culture & Society* 42, 3 (April 2020), 329–348. doi:10.1177/0163443719890530
- [114] Rodney T. Ogawa and Betty Malen. 1991. Towards Rigor in Reviews of Multivocal Literatures: Applying the Exploratory Case Study Method. *Review of Educational Research* 61, 3 (1991), 265–286. doi:10.3102/00346543061003265
- [115] Chinasa T. Okolo, Kehinde Aruleba, and George Obaido. 2023. *Responsible AI in Africa—Challenges and Opportunities*. Springer International Publishing, Cham, 35–64. doi:10.1007/978-3-031-08215-3\_3
- [116] Arnold Overwijk, Chenyan Xiong, Xiao Liu, Cameron VandenBerg, and Jamie Callan. 2022. ClueWeb22: 10 Billion Web Documents with Visual and Semantic Information. (2022). doi:10.48550/ARXIV.2211.15848
- [117] Ciaran O'Faircheallaigh. 2015. Social Equity and Large Mining Projects: Voluntary Industry Initiatives, Public Regulation and Community Development Agreements. *Journal of Business Ethics* 132, 1 (Nov. 2015), 91–103. doi:10.1007/s10551-014-2308-3
- [118] Arsenio Paez. 2017. Gray literature: An important resource in systematic reviews. *Journal of Evidence-Based Medicine* 10, 3 (2017), 233–240. doi:10.1111/jebm.12266
- [119] Matthew J. Page, Joanne E. McKenzie, Patrick M. Bossuyt, Isabelle Boutron, Tammy C. Hoffmann, Cynthia D. Mulrow, Larissa Shamseer, Jennifer M. Tetzlaff, Elie A. Akl, Sue E. Brennan, Roger Chou, Julie Glanville, Jeremy M. Grimshaw, Asbjørn Hróbjartsson, Manoj M. Lalu, Tianjing Li, Elizabeth W. Loder, Evan Mayo-Wilson, Steve McDonald, Luke A. McGuinness, Lesley A. Stewart, James Thomas, Andrea C. Tricco, Vivian A. Welch, Penny Whiting, and David Moher. 2021. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *BMJ* 372 (March 2021), n71. doi:10.1136/bmj.n71

- [120] Joon Sung Park, Danielle Bragg, Ece Kamar, and Meredith Ringel Morris. 2021. Designing an Online Infrastructure for Collecting AI Data From People With Disabilities. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Virtual Event Canada, 52–63. doi:10.1145/3442188.3445870
- [121] Yong Jin Park and S. Mo Jones-Jang. 2023. Surveillance, Security, and Ai as Technological Acceptance. *AI and Society* 38, 6 (2023), 2667–2678. doi:10.1007/s00146-021-01331-9
- [122] Frank Pasquale and Haochen Sun. 2024. Consent and Compensation: Resolving Generative AI’s Copyright Crisis. 4826695 (May 2024). doi:10.2139/ssrn.4826695
- [123] Michael Quinn Patton. 1991. Towards Utility in Reviews of Multivocal Literatures. *Review of Educational Research* 61, 3 (sept 1991), 287–292. doi:10.3102/00346543061003287
- [124] Charlotte Paul and Barbara Brookes. 2015. The Rationalization of Unethical Research: Revisionist Accounts of the Tuskegee Syphilis Study and the New Zealand “Unfortunate Experiment”. *American Journal of Public Health* 105, 10 (Oct. 2015), e12–19. doi:10.2105/AJPH.2015.302720
- [125] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only. arXiv:2306.01116 (2023). doi:10.48550/arXiv.2306.01116 arXiv:2306.01116 [cs].
- [126] Maneesha Perera, Rajith Vidanarachchi, Sangeetha Chandrashekeran, Melissa Kennedy, Brendan Kennedy, and Saman Halgamuge. 2025. Indigenous peoples and artificial intelligence: A systematic review and future directions. *Big Data & Society* 12, 2 (2025), 20539517251349170. doi:10.1177/20539517251349170
- [127] Claudio Santos Pinhanez and Edem Wornyo. 2025. Ethical Co-Development of AI Applications with Indigenous Communities. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA ’25)*. Association for Computing Machinery, New York, NY, USA, 1–4. doi:10.1145/3706599.3706649
- [128] Ian Pool. 2016. *Colonialism’s and postcolonialism’s fellow traveller: the collection, use and misuse of data on indigenous people* (1st ed.). ANU Press. doi:10.22459/CAEPR38.11.2016.04
- [129] Rida Qadri, Michael Madaio, and Mary L. Gray. 2025. Confusing the Map for the Territory – Communications of the ACM. <https://cacm.acm.org/opinion/confusing-the-map-for-the-territory/>
- [130] Rida Qadri, Piotr Mirowski, and Remi Denton. 2025. AI and Non-Western Art Worlds: Reimagining Critical AI Futures through Artistic Inquiry and Situated Dialogue. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–17. doi:10.1145/3706598.3714049
- [131] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21, 1 (Jan. 2020), 140:5485–140:5551.
- [132] Jenalea Rajab, Anuoluwapo Aremu, Evelyn Asiko Chimoto, Dale Dunbar, Graham Morrissey, Fadel Thior, Luandrie Potgieter, Jessico Ojo, Atnafu Lambebo Tonja, Maushami Chetty, Wilhelmina NdapewaOnyothi Nekoto, Pelonomi Moiloa, Jade Abbott, Vukosi Marivate, and Benjamin Rosman. 2025. The Esethu Framework: Reimagining Sustainable Dataset Governance and Curation for Low-Resource Languages. arXiv:2502.15916 (2025). doi:10.48550/arXiv.2502.15916 arXiv:2502.15916 [cs].
- [133] Mohammad Rashidujjaman Rifat, Abdullah Hasan Safir, Sourav Saha, Jahedul Alam Junaed, Maryam Saleki, Mohammad Ruhul Amin, and Syed Ishtiaque Ahmed. 2024. Data, Annotation, and Meaning-Making: The Politics of Categorization in Annotating a Dataset of Faith-based Communal Violence. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 2148–2156. doi:10.1145/3630106.3659030
- [134] Katja Rogers, Teresa Hirzle, Sukran Karaosmanoglu, Paula Toledo Palomino, Ekaterina Durmanova, Seiji Isotani, and Lennart E. Nacke. 2024. An Umbrella Review of Reporting Quality in CHI Systematic Reviews: Guiding Questions and Best Practices for HCI. *ACM Trans. Comput.-Hum. Interact.* 31, 5 (Nov. 2024), 57:1–57:55. doi:10.1145/3685266
- [135] Caroline Running Wolf and Noelani Arista. 2020. *Indigenous Protocols in Action*. Indigenous Protocol and Artificial Intelligence Working Group and the Canadian Institute for Advanced Research, Honolulu, HI, 93–101. doi:10.11573/spectrum.library.concordia.ca.00986506
- [136] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. “Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–15. doi:10.1145/3411764.3445518
- [137] Advait Sarkar. 2023. Enough With “Human-AI Collaboration”. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*. ACM, Hamburg Germany, 1–8. doi:10.1145/3544549.3582735
- [138] Devansh Saxena, Ji-Youn Jung, Jodi Forlizzi, Kenneth Holstein, and John Zimmerman. 2025. AI Mismatches: Identifying Potential Algorithmic Harms Before AI Development. arXiv:2502.18682 (April 2025). doi:10.48550/arXiv.2502.18682 arXiv:2502.18682 [cs].
- [139] Morgan Klaus Scheuerman, Alex Hanna, and Remi Denton. 2021. Do Datasets Have Politics? Disciplinary Values in Computer Vision Dataset Development. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (Oct. 2021), 1–37. doi:10.1145/3476058
- [140] Morgan Klaus Scheuerman, Kandrea Wade, Caitlin Lustig, and Jed R. Brubaker. 2020. How We’ve Taught Algorithms to See Identity: Constructing Race and Gender in Image Databases for Facial Analysis. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW1 (May 2020), 58:1–58:35. doi:10.1145/3392866
- [141] Morgan Klaus Scheuerman, Katy Weathington, Tarun Mugunthan, Emily Denton, and Casey Fiesler. 2023. From Human to Data to Dataset: Mapping the Traceability of Human Subjects in Computer Vision Datasets. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW1 (April 2023), 1–33. doi:10.1145/3579488

- [142] Morgan Klaus Scheuerman, Allison Woodruff, and Jed R. Brubaker. 2025. How Data Workers Shape Datasets: The Role of Positionality in Data Collection and Annotation for Computer Vision. *Proc. ACM Hum.-Comput. Interact.* 9, 7 (Oct. 2025), CSCW300:1–CSCW300:42. doi:10.1145/3757481
- [143] Daniel Schiff, Bogdana Rakova, Aladdin Ayesh, Anat Fanti, and Michael Lennon. 2020. *Principles to Practices for Responsible AI: Closing the Gap*. doi:10.48550/arXiv.2006.04707 arXiv:2006.04707 [cs].
- [144] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. (2022). doi:10.48550/ARXIV.2210.08402
- [145] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. LAION-400M: Open Dataset of CLIP-Filtered 400 Million Image-Text Pairs. (2021). doi:10.48550/ARXIV.2111.02114
- [146] Candice Schumann, Gbolahan O. Olanubi, Auriel Wright, Ellis Monk Jr., Courtney Heldreth, and Susanna Ricco. 2023. Consensus and Subjectivity of Skin Tone Annotation for ML Fairness. <https://arxiv.org/abs/2305.09073v3>
- [147] Eve Kosofsky Sedgwick. 1990. *Epistemology of the Closet*. University of California Press. Google-Books-ID: u5jgaOhhmpgC.
- [148] Judith Shapiro and John-Andrew McNeish. 2021. *Our Extractive Age: Expressions of Violence and Resistance* (1 ed.). Routledge, London. doi:10.4324/9781003127611
- [149] Renee Shelby, Shalaleh Rismani, Kathryn Henne, AJung Moon, Negar Rostamzadeh, Paul Nicholas, N’Mah Yilla, Jess Gallegos, Andrew Smart, Emilio Garcia, and Gurleen Virk. 2023. Sociotechnical Harms of Algorithmic Systems: Scoping a Taxonomy for Harm Reduction. arXiv:2210.05791 (2023). <http://arxiv.org/abs/2210.05791> arXiv:2210.05791 [cs].
- [150] Jan Simson, Alessandro Fabris, and Christoph Kern. 2024. Lazy Data Practices Harm Fairness Research. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 642–659. doi:10.1145/3630106.3658931
- [151] Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Devidas Mataciunas, Laura OMahony, et al. 2024. Aya dataset: An open-access collection for multilingual instruction tuning. *arXiv preprint arXiv:2402.06619* (2024).
- [152] Linda Tuhiwai Smith. 2021. *Decolonizing Methodologies: Research and Indigenous Peoples* (third edition ed.). Bloomsbury Publishing. Google-Books-ID: EwA1EAAQBAJ.
- [153] Gayatri Chakravorty Spivak. 1994. *Can the Subaltern Speak?* Routledge, London, 66–111.
- [154] Yolande Strengers, Jathan Sadowski, Zhuying Li, Anna Shimshak, and Florian “Floyd” Mueller. 2021. What Can HCI Learn from Sexual Consent?: A Feminist Process of Embodied Consent for Interactions with Emerging Technologies. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–13. doi:10.1145/3411764.3445107
- [155] Harini Suresh and John V. Gutttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In *Equity and Access in Algorithms, Mechanisms, and Optimization*. 1–9. doi:10.1145/3465416.3483305 arXiv:1901.10002 [cs, stat].
- [156] Harini Suresh, Emily Tseng, Meg Young, Mary Gray, Emma Pierson, and Karen Levy. 2024. Participation in the age of foundation models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1609–1621.
- [157] Alex S. Taylor. 2011. Out there. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI ’11)*. Association for Computing Machinery, New York, NY, USA, 685–694. doi:10.1145/1978942.1979042
- [158] Jordan Taylor, Wesley Hanwen Deng, Kenneth Holstein, Sarah Fox, and Haiyi Zhu. 2024. Carefully Unmaking the “Marginalized User:” A Diffraction Analysis of a Gay Online Community. *ACM Transactions on Computer-Human Interaction* (2024), 3673229. doi:10.1145/3673229
- [159] John Taylor and Tahu Kukutai (Eds.). 2016. *Indigenous Data Sovereignty: Toward an Agenda*. ANU Press, Acton, ACT, Australia. doi:10.22459/CAEPR38.11.2016
- [160] Te Hiku Media. [n. d.]. *Kaitiakitanga License*. <https://github.com/TeHikuMedia/Kaitiakitanga-License> GitHub repository.
- [161] Jaime A. Teixeira da Silva. 2022. Handling Ethics Dumping and Neo-Colonial Research: From the Laboratory to the Academic Literature. *Journal of Bioethical Inquiry* 19, 3 (2022), 433–443. doi:10.1007/s11673-022-10191-x
- [162] Tesfa-Alem Tekle. 2020. Refugees in Ethiopia’s camps raise privacy and exclusion concerns over UNHCR’s new digital registration. <https://globalvoices.org/2020/03/19/refugees-in-ethiopia-camps-raise-privacy-and-exclusion-concerns-over-unhcrs-new-digital-registration/>
- [163] Jim Thatcher, David O’Sullivan, and Dillon Mahmoudi. 2016. Data colonialism through accumulation by dispossession: New metaphors for daily data. *Environment and Planning D: Society and Space* 34, 6 (Dec. 2016), 990–1006. doi:10.1177/0263775816633195
- [164] Scott Timcke. 2024. AI and the digital scramble for Africa. <https://roape.net/2024/07/11/ai-and-the-digital-scramble-for-africa/>
- [165] Atnafu Lambebo Tonja, Bonaventure F. P. Dossou, Jessica Ojo, Jenalea Rajab, Fadel Thior, Eric Peter Wairagala, Anuoluwapo Aremu, Pelonomi Moilola, Jade Abbott, Vukosi Marivate, and Benjamin Rosman. 2024. InkubaLM: A small language model for low-resource African languages. arXiv:2408.17024 (2024). doi:10.48550/arXiv.2408.17024 arXiv:2408.17024 [cs].
- [166] Paola Tubaro, Antonio A. Casilli, Maxime Cornet, Clément Le Ludec, and Juana Torres Cierpe. 2025. Where does AI come from? A global case study across Europe, Africa, and Latin America. (Feb. 2025). doi:10.1080/13563467.2025.2462137 arXiv:2502.04860 [cs].
- [167] Ding Wang, Shantanu Prabhat, and Nithya Sambasivan. 2022. Whose AI Dream? In search of the aspiration in data annotation. arXiv:2203.10748 (March 2022). doi:10.48550/arXiv.2203.10748 arXiv:2203.10748 [cs].
- [168] Lining Wang, Vaishnav Kameswaran, and Hernisa Kacorri. 2025. Toward a Taxonomy of Algorithmic Harms for Disability: A Systematic Review. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* 8, 3 (Oct. 2025), 2649–2665. doi:10.1609/aies.v8i3.36745

- [169] Karl Werder, Balasubramaniam Ramesh, and Rongen (Sophia) Zhang. 2022. Establishing Data Provenance for Responsible Artificial Intelligence Systems. *ACM Transactions on Management Information Systems* 13, 2 (2022), 1–23. doi:10.1145/3503488
- [170] Cedric Deslandes Whitney and Justin Norman. 2024. Real Risks of Fake Data: Synthetic Data, Diversity-Washing and Consent Circumvention. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 1733–1744. doi:10.1145/3630106.3659002
- [171] David Gray Widder. 2024. Epistemic Power in AI Ethics Labor: Legitimizing Located Complaints. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. ACM, Rio de Janeiro Brazil, 1295–1304. doi:10.1145/3630106.3658973
- [172] David Gray Widder, Meredith Whittaker, and Sarah Myers West. 2024. Why ‘open’ AI systems are actually closed, and why this matters. *Nature* 635, 8040 (Nov. 2024), 827–833. doi:10.1038/s41586-024-08141-1
- [173] Claes Wohlin. 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In *Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering*. ACM, London England United Kingdom, 1–10. doi:10.1145/2601248.2601268
- [174] Lucy C. Woodall, Sheena Talma, Oliver Steeds, Paris Stefanoudis, Marie-May Jeremie-Muzungaile, and Alain de Comarmond. 2021. Co-development, co-production and co-dissemination of scientific research: a case study to demonstrate mutual benefits. *Biology Letters* 17, 4 (April 2021), 20200699. doi:10.1098/rsbl.2020.0699
- [175] Meg Young, Upol Ehsan, Ranjit Singh, Emnet Tafesse, Michele Gilman, Christina Harrington, and Jacob Metcalf. 2024. Participation versus scale: Tensions in the practical demands on participatory AI. *First Monday* (April 2024). doi:10.5210/fm.v29i4.13642
- [176] Edibe Betul Yucer. 2025. AI is finally trying to speak African languages. Will this end a historic neglect? *TRT Global* (Aug. 2025). <https://trt.global/afrika-english/article/359e1362af39>
- [177] Dora Zhao, Morgan Klaus Scheuerman, Pooja Chitre, Jerone T. A. Andrews, Georgia Panagiotidou, Shawn Walker, Kathleen H. Pine, and Alice Xiang. 2024. A taxonomy of challenges to curating fair datasets. In *Proceedings of the 38th International Conference on Neural Information Processing Systems (NIPS ’24, Vol. 37)*. Curran Associates Inc., Red Hook, NY, USA, 97826–97858.

## A Review Methodologies in HCI

Review methodologies in HCI face persistent challenges when systematic transparency must coexist with interpretive expertise and when evidence circulates across fragmented publication ecosystems [48, 134]. Systematic reviews establish protocols for organizing evidence within bounded domains, and PRISMA frameworks support reproducibility [119]. Both approaches reach limits in interdisciplinary settings where evidence types vary and epistemological frameworks conflict [6, 109].

Multivocal literature reviews (MLRs) offer one way to meet these demands and balance rigor with practical utility, as recent work in responsible AI demonstrates [95]. The approach originated in educational research as a methodological framework to impose systematic rigor on reviews of diverse documents [114], initiating a discussion that clarified that the standard of rigor must be situational and secondary to utility for practitioners [123]. Software engineering later adapted MLRs to capture both state-of-the-art research and state-of-practice knowledge [52, 53]. MLRs integrate peer-reviewed academic literature with grey literature such as organizational reports, policy documents, community statements, technical documentation, practitioner outputs, and multimedia materials. The widely cited Luxembourg definition characterizes grey literature as material produced by government, academia, industry, or community groups that is not controlled by commercial publishers [55]. Diversity of source types and timeliness are core advantages, since emerging practices often circulate outside formal publication channels and appear earlier than peer-reviewed work [118].

Credibility varies across grey-literature types, and assessments often depend on provenance, expertise, and recognized authority [74]. Multivocal approaches are particularly important for scholarship involving Indigenous and underserved communities. Lewis et al. [89] show that multivocality preserves heterogeneous viewpoints. The authors combine essays, protocols, and artistic works rather than imposing a single scholarly mode, a stance aligning with decolonizing methodologies that emphasize community-generated knowledge and Indigenous epistemic authority [10, 152]. Additionally, influential contributions in AI ethics and data governance often appear in organizational reports [25, 69], investigative journalism [63], public initiatives [1], and widely cited preprints [15]. MLR methods accommodate a diversity of distributed knowledge production and support synthesis across venues not fully captured by academic indexing.

## B Search Strategy & Corpus Composition

Database searches were conducted iteratively between August 2024 and January 2025, complementing network referrals and citation snowballing. The final structured ACM Digital Library search was executed on January 31, 2025, using the advanced search interface with abstract and full-text indexing via personal subscription. Table 3 reports the four primary ACM query sets and their outcomes.

Aggregate results from the January 2025 ACM searches are summarized in Table 5. Across 1,914 hits, 1,201 items were screened, yielding 153 that met criteria and 48 unique sources after full-text review and duplicate removal. Similar comprehensive search strategies were applied to IEEE Xplore, ScienceDirect, Taylor & Francis Online, Wiley Online Library, Google Scholar, and Springer Link, following the same phased approach for queries.

Table 3. ACM Digital Library search queries and results.

| Query | Exact search string  | Results screened              | Potentially relevant            |
|-------|--|-------------------------------|---------------------------------|
| Q1    | ((("data collection" OR "data production" OR "data curation" OR "dataset development") AND ("artificial intelligence" OR "machine learning" OR "AI") AND ("marginalized" OR "underrepresented" OR "underserved" OR "community" OR "indigenous")) | 51 abstracts + 300 full-texts | 45 (7 abstracts, 38 full-texts) |
| Q2    | ((("extractive" OR "exploitative" OR "data colonialism") AND ("data practices" OR "dataset construction") AND ("communities" OR "workers" OR "labor"))   | 3 abstracts + 182 full-texts  | 13 (1 abstract, 12 full-texts)  |
| Q3    | ((("crowdsourcing" OR "platform labor") AND ("bias" OR "fairness" OR "ethics") AND ("marginalized" OR "vulnerable populations" OR "community harm"))   | 491 full-texts (200 screened) | 32                              |
| Q4    | ((("participatory design" OR "community-led" OR "co-design") AND ("ai development" OR "dataset creation") AND ("sovereignty" OR "community engagement" OR "ethical data"))   | 122 full-texts                | 12                              |

Table 4. Targeted ACM venue-specific searches.

| Venue                      | Exact search string  | Venue filter   | Results summary   |
|----------------------------|--|--|---|
| CHI Conference Proceedings | ((("data collection" OR "data production" OR "data curation" OR "dataset development") AND ("artificial intelligence" OR "machine learning" OR "AI") AND ("marginalized" OR "underrepresented" OR "underserved" OR "community" OR "indigenous")) | CHI Conference on Human Factors in Computing Systems (all years) | 296 hits; screened: first 200; potentially relevant: 15 |
| FAccT Proceedings          | ((("data collection" OR "data production" OR "data curation" OR "dataset development") AND ("artificial intelligence" OR "machine learning" OR "AI") AND ("marginalized" OR "underrepresented" OR "underserved" OR "community" OR "indigenous")) | ACM Conference on Fairness, Accountability, and Transparency     | 143 hits; screened: all; potentially relevant: 37       |

Table 5. ACM search results summary.

| Category  | Count |
|---|-------|
| Total primary searches (query sets)                       | 6     |
| Total venue-specific searches                             | 2     |
| Total hits across all searches                            | 1,914 |
| Total items screened (varied by search size)              | 1,201 |
| Items meeting inclusion criteria after screening          | 153   |
| Items retained after full-text review                     | 89    |
| Final unique sources for corpus (after duplicate removal) | 48    |

Table 6. Discovery method distribution (N=350 sources).

| Method                          | Sources    | Percent     |
|---------------------------------|------------|-------------|
| Database searches               | 174        | 50%         |
| Existing networks/organizations | 73         | 21%         |
| Citation snowballing            | 51         | 15%         |
| Iterative keyword search        | 31         | 9%          |
| Hand-searching journals         | 21         | 6%          |
| <b>Total</b>                    | <b>350</b> | <b>100%</b> |

## C Corpus Creation Details

### Datasheet Fields.

Table 7. Description of datasheet fields

| Column                 | Content  |
|------------------------|--|
| <b>Identifier</b>      | In-line APA citation (author surname and year) used as a unique ID for tracking within the corpus.   |
| <b>APA Citation</b>    | Full APA reference for the source.   |
| <b>Title</b>           | Title of the publication or output.  |
| <b>Analytic Domain</b> | Controlled list, multiple possible: <ul style="list-style-type: none"> <li>• Data Relations: how data is scoped, justified, and negotiated</li> <li>• Data Labor: how curation work is arranged and carried out</li> <li>• Data Representation: how categories are constructed and based on what presences and absences</li> <li>• Data Infrastructure: how technical systems mediate data movement</li> <li>• Data Governance: how authority over data shapes downstream use</li> </ul> |
| <b>Orientation</b>     | Controlled list: <ul style="list-style-type: none"> <li>• Extractive: undermines consent, compensation, or benefit</li> <li>• High-Agency Principles: normative frameworks promoting stewardship, sovereignty, accountability</li> <li>• High-Agency Practices: operationalized, community-led, participatory, or sovereignty-based initiatives</li> </ul>   |
| <b>General Theme</b>   | Controlled list, multiple possible: <ul style="list-style-type: none"> <li>• Community impacts and relations</li> <li>• Critical theory</li> <li>• Data labor</li> <li>• Data practices</li> <li>• Ethics frameworks</li> </ul>  |

| Column                                 | Content  |
|--|--|
| <b>Pipeline Stage</b>                  | Specific process within a pipeline stage (controlled list): <ul style="list-style-type: none"> <li>• Problem Understanding and Formulation</li> <li>• ML System Design and Development</li> <li>• Deployment and Impact</li> <li>• Cross-pipeline</li> </ul>   |
| <b>Pipeline Sub-stage</b>              | Specific process within a pipeline stage (controlled list): <ul style="list-style-type: none"> <li>• Institutional Prioritization and Funding</li> <li>• Product Conception and Design</li> <li>• Data Selection, Collection and Annotation</li> <li>• Model Architecture Selection and Design</li> <li>• Model Training and Evaluation</li> <li>• Product Testing</li> <li>• Product Launch</li> <li>• Cross-pipeline</li> </ul>  |
| <b>Historical Era</b>                  | Era of data production practice (controlled list): <ul style="list-style-type: none"> <li>• Era 1: Curated datasets (pre-2009); no sources in corpus</li> <li>• Era 2: Crowdsourced benchmarks (2009–2017)</li> <li>• Era 3: Web-scraped/foundation models (2017–present)</li> <li>• Multi-era: Spans multiple eras or provides historical analysis</li> </ul>   |
| <b>Primary Pattern(s) / Pathway(s)</b> | The specific extractive or high-agency behavior described in the source. Between one and three tags were assigned per source in order of relevance. For sources that provide conceptual, historical, or framing contributions without mapping directly onto an identified pattern, we assigned Other/NA (conceptual framing).  |
| <b>Triangle Coverage</b>               | Engagement with the three scoping domains that defined corpus eligibility: <ul style="list-style-type: none"> <li>• A — AI contexts</li> <li>• D — Data production practices</li> <li>• C — Community impacts</li> </ul> <p>Because community impacts (C) establish the outer bounds of the review, included sources substantively address all three domains, though with varying emphases. Codes (A, D, C or combinations ADC, DC, AD) indicate which domains are explicitly developed within the source. An accompanying Rationale column explains the basis for inclusion and the specific ways each source engages A/D/C beyond passing mention.</p> |

| Column                                      | Content   |
|---|---|
| <b>How Source Was Found</b>                 | White literature (journal papers, conference proceedings, books) or Grey literature (reports, policy documents, theses, community outputs, blogs).  |
| <b>Keywords</b>                             | 3–5 terms for coding/search, ordered Geography → Data/technical → Community/impact.   |
| <b>Geographic Region of Focus</b>           | Region or community under study (controlled list): Africa, APAC, EU/UK, LatAm, MENA, North America, Oceania, Multiple regions, Not regionally specific (globally framed advocacy, transnational collectives, or technical works not tied to one region).          |
| <b>Author Affiliations</b>                  | High-level institutional grouping of authors. If multiple affiliations, code majority grouping here; record full details in Authorship & Positionality Context. Controlled list: Academic; Government; Industry; NGO/Non-profit; Mixed; Journalist/Other/Not sure |
| <b>Geographic Area of Author(s)</b>         | Full institution name and country of the lead author(s)   |
| <b>Institution</b>                          | Region of lead author's institution (controlled list, same regions as above).   |
| <b>Authorship and Positionality Context</b> | Complete authorship profile, including all institutions, geographic distribution, equal contribution notes, and any relevant statements on positionality or disciplinary traditions.  |
| <b>Summary</b>                              | ≤ 120-word synopsis. Structure: Topic → Method → Findings → Link to AI data production + community impacts.   |

**Corpus Summary.**

Table 8. Corpus composition summary (N=350 sources)

| Category                       | Sub-category                        | Count (%) |
|--------------------------------|-------------------------------------|-----------|
| Orientation                    | Extractive Practices                | 141 (41%) |
|                                | High-Agency Principles              | 116 (33%) |
|                                | High-Agency Practices               | 93 (27%)  |
| Source Type                    | White literature                    | 258 (74%) |
|                                | Grey literature                     | 92 (26%)  |
| Geographic Focus               | Not regionally specific             | 150 (43%) |
|                                | Multiple areas                      | 61 (17%)  |
|                                | North America                       | 41 (11%)  |
|                                | Africa                              | 38 (11%)  |
|                                | Oceania                             | 19 (5%)   |
|                                | APAC                                | 15 (4%)   |
|                                | EU/UK                               | 11 (3%)   |
|                                | LatAm                               | 12 (3%)   |
|                                | MENA                                | 3 (1%)    |
| Author Affiliation (lead only) | Academic                            | 182 (52%) |
|                                | Mixed                               | 96 (27%)  |
|                                | Industry                            | 37 (11%)  |
|                                | NGO/Non-profit                      | 20 (6%)   |
|                                | Journalist/Other                    | 11 (3%)   |
|                                | Government                          | 4 (1%)    |
| Pipeline Stage                 | ML System Design & Development      | 183 (52%) |
|                                | Problem Understanding & Formulation | 90 (25%)  |
|                                | Cross-pipeline                      | 55 (16%)  |
|                                | Deployment & Impact                 | 22 (6%)   |
| Historical Era                 | Era 3 (2017–present)                | 241 (69%) |
|                                | Multi-era                           | 95 (27%)  |
|                                | Era 2 (2009–2017)                   | 14 (4%)   |
|                                | Era 1 (pre-2009)                    | 0 (0%)    |

Received 4 September 2025; revised 4 December 2025; accepted 5 June 2009