



Expanding participatory governance for LLMs

Case studies from BigCode, Cohere for AI's Aya Initiative, and Collective Intelligence Project

Jennifer Ding | The Alan Turing Institute

March 2024



Tools, practices and systems

Building open source infrastructure to empower a global, decentralised network of people who connect data with domain experts

Learn more ↓

Introduction

Data science and artificial intelligence are becoming ever more prevalent in the UK with common needs and challenges arising. Solving these challenges requires novel tools, practices and systems which can unlock advances across the wider sector and accelerate innovation.

The tools, practices and systems (TPS) programme at the Turing represents a cross-cutting set of initiatives which seek to build open source infrastructure that is accessible to all, and to empower a global, decentralised network of people who connect data with domain experts.

The programme will:

- Build trustworthy systems
- Embed transparent reporting practices
- Promote inclusive interoperable design
- Maintain ethical integrity
- Encourage respectful co-creation

Community Management

Research Application Management

Participatory Science

Organisers

Contact info

Trustworthy Systems

The Turing Way

Data Safe Haven

Pitchfest

London Data Week

Urban Analytics Tech Platform

AIM-RSF

AutSPACES

Turing Commons

Trustworthy Ethical Assurance (TEA)

TRIC Impact Hub

Agenda

1. Motivations
2. Case Studies: BigCode, Aya, CIP
3. Current Work & Emerging Projects

Agenda

1. *Motivations*
2. Case Studies: BigCode, Aya, CIP
3. Current Work & Emerging Projects

**OPENCV +
SCIKIT-LEARN**

BERT

YOLO

OPEN LLMs

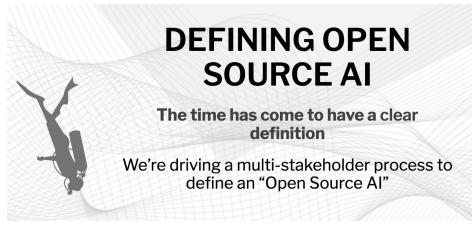


**The
Alan Turing
Institute**



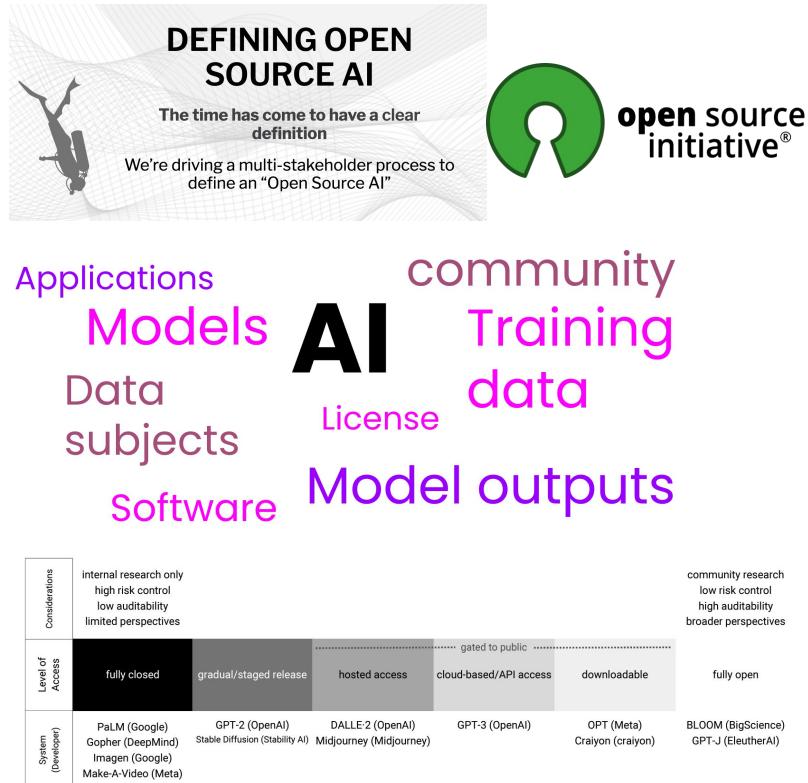
What defines the ‘Open’ in ‘Open Source AI’?

- Stable descriptor as a shared foundation for tools, regulation, best practices



What defines the ‘Open’ in ‘Open Source AI’?

- Stable descriptor as a shared foundation for tools, regulation, best practices
- Designing rather than retrofitting “open” and “open source” for AI
 - Thinking beyond models
 - Interactions with “responsible” and “safe” AI
 - Open as a gradient rather than a binary



Who is building open source AI?

21 September, 14:00 - 15:30 UTC+1

Register on Eventbrite



Arielle Bennett

Programme Manager
The Alan Turing Institute



Mophat Okinyi

Union Representative
African Content Moderators Union



Marzieh Fadaee

Senior Research Scientist
Cohere for AI



Abinaya Mahendiran

CTO
Nunnari Labs



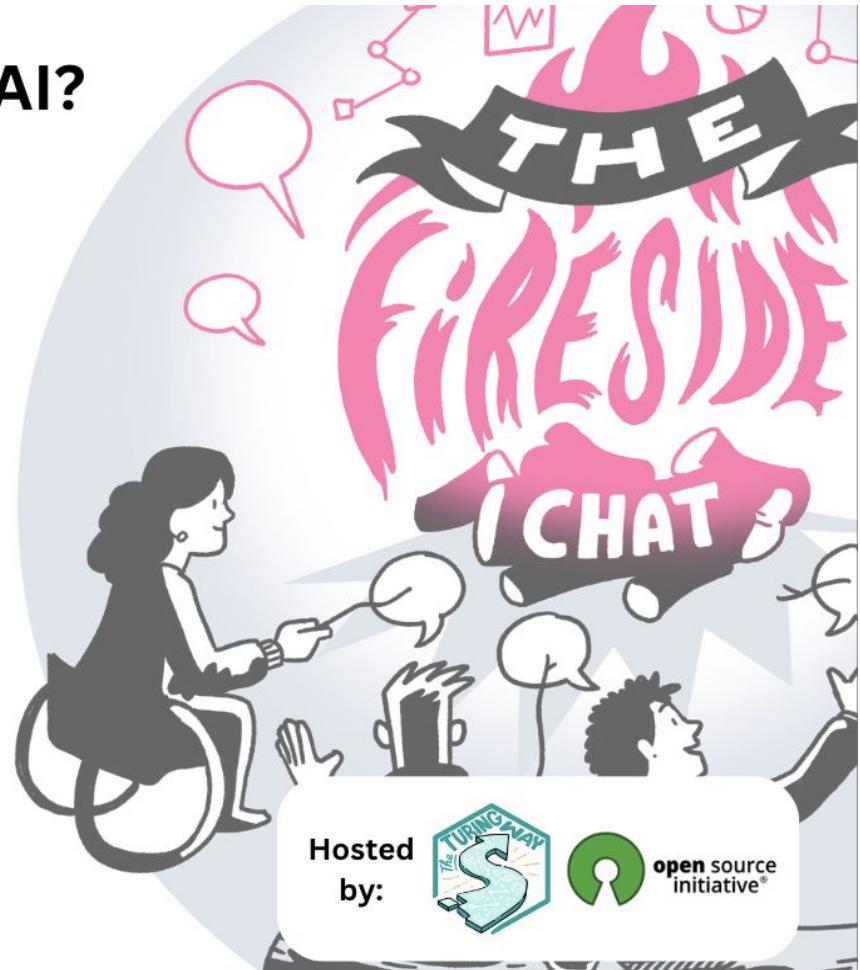
David Gray Widder

Postdoctoral Fellow
Cornell Tech



Jennifer Ding

Senior Researcher
The Alan Turing Institute



Open Source AI Collaborative

- Distributed, volunteer-led teams creating alternative AI development pathways grounded in “open” principles
- Focus on different activities relevant for their respective communities:
 - building code LLMs (**BigCode**)
 - building multilingual LLMs (BigScience Workshop, **Aya Initiative**)
 - creating model constitutions (**Collective Intelligence Project**)

Towards Openness Beyond Open Access: User Journeys through 3 Open AI Collaboratives

Jennifer Ding
The Alan Turing Institute
jing@turing.ac.uk

Christopher Akiki
Leipzig University
christopher.akiki@uni-leipzig.de

Yacine Jernite
Hugging Face
yacine@huggingface.co

Anne Lee Steele
The Alan Turing Institute
asteelle@turing.ac.uk

Temi Popo
Mozilla Foundation
temi@mozilla.org

Abstract

Open Artificial Intelligence (Open AI) collaboratives offer alternative pathways for how AI can be developed beyond well-resourced technology companies and who can be a part of the process. To understand how and why they work and what added value they bring to the field, we conducted three case studies, each focused on a different kind of activity around AI building code (BigScience workshop), *methodways of working* (The Turing Way), and *ecosystems* (Mozilla Festival's Building Trustworthy AI Working Group). First, we document the community structures that facilitate these distributed, volunteer-led teams, comparing the collaboration styles that drive each group towards their specific goals. Through interviews with community leaders, we map user journeys for how members discover, join, contribute, and participate. Ultimately, this paper aims to highlight the diversity of AI work and workers and how it contributes to these collaborations and how they offer a broader practice of openness to the AI space.



← Cohere For AI |

Languages are not treated equally by researchers. Some languages have received disproportionate attention and focus in NLP.

Language	# of papers per million speakers	# of speakers (in millions)
Irish	5235	0.2
Basque	2430	0.5
German	179	83
English	63	550
Chinese	11	1000
Hausa	1.5	70
Nigerian Pidgin	0.4	30

Van Esch et al. 2022

Agenda

1. Motivations
2. ***Case Studies: BigCode, Aya, CIP***
3. Current Work & Emerging Projects

Case Study: BigCode & StarCoder LLM

“Open not only for transparency but accountability”



The Stack is an open governance interface between the AI community and the open source community.

Am I in The Stack?

As part of the BigCode project, we released and maintain [The Stack](#), a 6 TB dataset of permissively licensed source code over 300 programming languages. One of our goals in this project is to give people agency over their source code by letting them decide whether or not it should be used to develop and evaluate machine learning models, as we acknowledge that not all developers may wish to have their data used for that purpose.

This tool lets you check if a repository under a given username is part of The Stack dataset. Would you like to have your data removed from future versions of The Stack? You can opt-out following the instructions [here](#).

The Stack version:

v1.2

Your GitHubs username:

dingaling

Yes, there is code from 6 repositories in The Stack:

[dingaling/citi-map](#)

[dingaling/simple-ChatBot](#)

[dingaling/street-light-map](#)

[dingaling/the-turing-way](#)

[dingaling/tramTracker](#)

[dingaling/turing-commons](#)

Opt-out

If you want your data to be removed from the stack and model training open an issue with [this link](#) (if the link doesn't work try right a right click and open it in a new tab) or visit <https://github.com/bigcode-project/opt-out-v2/issues/new>

[&template=opt-out-request.md](#)

bigcode/starcoder 2.42k

Text Generation Transformers PyTorch bigcode/the-stack-dedup gpt_bigcode code Eval Results Inference Endpoints text-generation-inference arxiv:1911.02150 arxiv:2205.14135

anvil:2207.14255 anvil:2305.06161 License: bigcode-openrail-m

Model card Files and versions Community

You need to agree to share your contact information to access this model

This repository is publicly accessible, but you have to accept the conditions to access its files and content.

Model License Agreement

Please read the BigCode [OpenRAILM license](#) agreement before accepting it.

Log in or Sign Up to review the conditions and access this model content.

StarCoder

StarCoder

How to join?

We are excited to invite AI practitioners from diverse backgrounds to join the BigCode project! Note that BigCode is a *research collaboration* and is open to participants who

1. have a professional research background and
2. are able to commit time to the project.

In general, we expect applicants to be affiliated with a research organization (either in academia or industry) and work on the technical/ethical/legal aspects of LLMs for coding applications.

You can apply here to the BigCode project!

BigCode Project Governance Card

THE BIGCODE PROJECT GOVERNANCE CARD

Sean Hughes^{1,*} Harm de Vries² Jennifer Robinson¹
Carlos Muñoz Ferrandis^{3,*} Loubna Ben Allal³ Leandro von Werra³
Jennifer Ding⁴ Sébastien Paquet² Yacine Jernite³

¹ServiceNow ²ServiceNow Research ³Hugging Face ⁴The Alan Turing Institute

Corresponding authors (*) can be contacted at contact@bigcode-project.org

Abstract

This document serves as an overview of the different mechanisms and areas of governance in the BigCode project. It aims to support transparency by providing relevant information about choices that were made during the project to the broader public, and to serve as an example of intentional governance of an open research project that future endeavors can leverage to shape their own approach. The first section, Project Structure, covers the project organization, its stated goals and values, its internal decision processes, and its funding and resources. The second section, Data and Model Governance, covers decisions relating to the questions of data subject consent, privacy, and model release.



The Turing Way

Search this book...

Welcome
Guide for Reproducible Research
Guide for Project Design
Overview of Project Design
Project Design Checklist
Creating Project Repositories
Personas and Pathways
File Naming Convention
Code Styling and Linting
Sensitive Data Projects
Managing Sensitive Data
Projects
Working on Sensitive Data
Projects
Data Governance
Data Governance for the Machine Learning Pipeline
BigCode Data Governance
Case Study





































































































































































































































































































































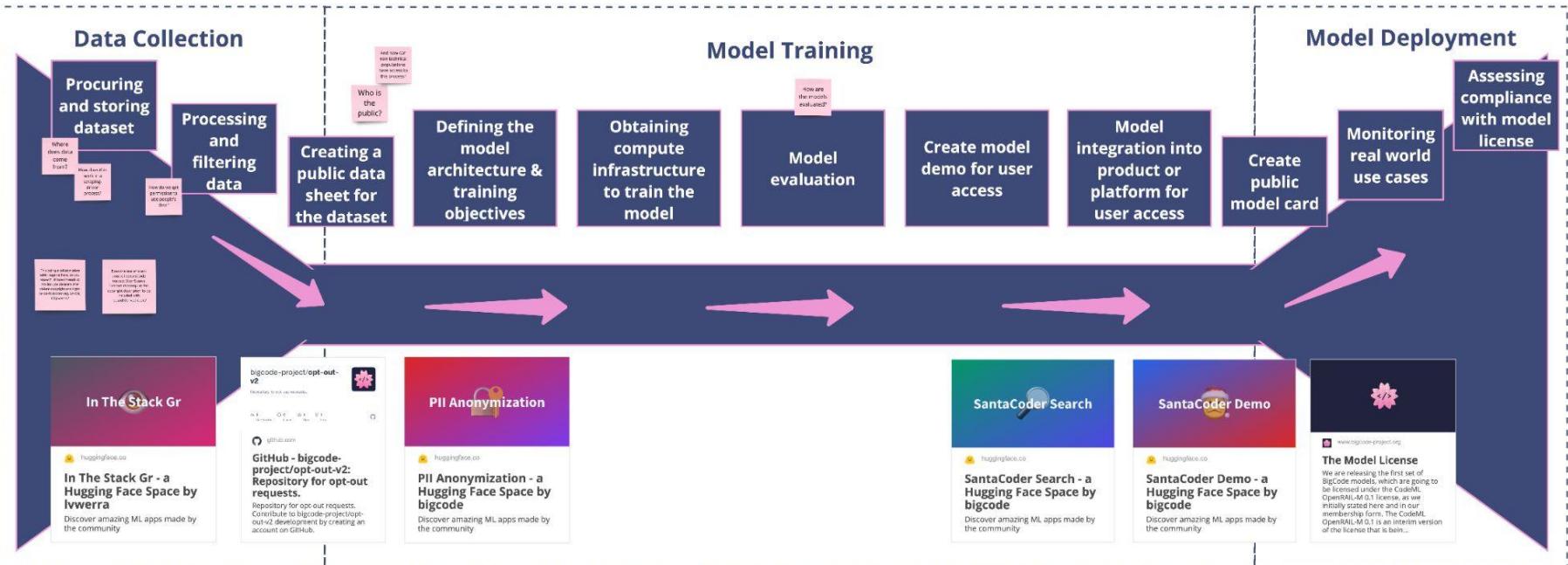





StarCoder LLM Governance Pipeline

“Open not only for transparency but accountability”

case study: BigCode



Case Study: Cohere for AI's Aya Initiative

“Harness the collective wisdom and contributions of people from all over the world.”



(a) Example of an original annotation contribution.

(b) Example of a re-annotation contribution.

Aya Initiative Data Paper with Language Ambassadors

The screenshot shows the Cohere website with a blue header bar. On the left is the Cohere logo. To its right are five navigation links: PRODUCTS, FOR BUSINESS, DEVELOPERS, RESEARCH, and COMPANY. Below the header, there's a back-to-more-papers link and the title "Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning". A "MULTILINGUAL" tag is present. At the bottom, there's a "READ THE PAPER" button with a right-pointing arrow icon.



Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning

Shivalika Singh^{*1}, Freddie Vargas^{*1}, Daniel D'souza^{*1}, Börje F. Karlsson^{*2},
Abinaya Mahendiran^{*1}, Wei-Yin Ko^{*3}, Herumb Shandilya^{*4}, Jay Patel^{*4},
Deividas Mataciunas^{*5}, Laura O'Mahony^{*5}, Mike Zhang^{*6}, Ramith Hettiarachchi^{*7},
Joseph Wilson^{*8}, Marina Machado^{*9}, Lissa Souza Moura^{*9}, Dominik Krzeminski^{*9},
Hakimeld Fadace^{*10}, Irem Ergun^{*11}, Heoom Okoh^{*11}, Aisha Alaagib^{*12},
Oshan Mudannayake^{*13}, Alyafeai^{*13}, Maha Chhetri^{*14}, Sebastian Ruder^{*15},
Surya Guthikonda^{*16}, Eman Alghamdi^{*16}, Schuyler Gehrmann^{*11},
Niklas Muennighoff^{*17}, Max Bartolo^{*18}, Julia Kreutzer^{*12}, Ahmed Ustun^{*12},
Marzich Fadace^{*12}, and Sara Hooker^{*12}

^{*1}Cohere For AI Community, ^{*2}Beijing Academy of Artificial Intelligence, ^{*3}Cohere, ^{*4}Binghamton University,
^{*5}University of Limerick, ^{*6}TU University of Copenhagen, ^{*7}MIT, ^{*8}University of Toronto, ^{*9}King Fahd University of
Petroleum and Minerals, ^{*10}King Abdulaziz University, ASAS.AE, ^{*11}Bloomberg LP, ^{*12}Cohere For AI

Corresponding authors: Shivalika Singh <sivalikasingh96@gmail.com>, Marzich Fadace <marzich@cohere.com>,
Sara Hooker <sarahooker@cohere.com>

Abstract

Datasets are foundational to many breakthroughs in modern artificial intelligence. Many recent achievements in the space of natural language processing (NLP) can be attributed to the fine-tuning of pre-trained models on a diverse set of tasks that enables a large language model (LLM) to respond to instructions. Instruction fine-tuning (IFT) requires specifically constructed and annotated datasets. However, existing datasets are almost all in the English language. In this work, our primary goal is to bridge the language gap by building a human-curated instruction-following dataset spanning 65 languages. We worked with fluent speakers of languages from around the world to collect natural instances of instructions and completions. Furthermore, we create the most extensive multilingual collection to date, comprising 513 million instances through templating and translating existing datasets across 114 languages. In total, we contribute four key resources: we develop and open-source the [Aya Annotation Platform](#), the [Aya Dataset](#), the [Aya Collection](#), and the [Aya Evaluation Suite](#). The Aya initiative also serves as a valuable case study in participatory research collaborations involving collaborators from 119 countries. We see this as a valuable framework for future research collaborations that aim to bridge gaps in resources.

Singh, S., Vargas, F., Dsouza, D., Karlsson, B.F., Mahendiran, A., Ko, W., Shandilya, H., Patel, J., Mataciunas, D., OMahony, L., Zhang, M., Hettiarachchi, R., Wilson, J., Machado, M., Moura, L.S., Krzeminski, D., Fadaei, H., Ergun, I., Okoh, I., Alaagib, A., Mudannayake, O., Alyafeai, Z., Vu, M.C., Ruder, S., Guthikonda, S., Alghamdi, E.A., Gehrmann, S., Muennighoff, N., Bartolo, M., Kreutzer, J., Ustun, A., Fadaee, M., & Hooker, S. (2024). Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning. ArXiv, abs/2402.06619.

Case Study: Collective Intelligence Project

Effective, decentralized, and agentic decision-making across individuals and communities to produce best-case outcomes for the collective.

The Collective Intelligence Project

Whitepaper Research Blog About 

Table of Contents

[Introducing the Collective Intelligence Project](#)

[The Transformative Technology Trilemma](#)

- I. Capitalist Acceleration: Sacrificing safety for progress
- II. Authoritarian Technocracy: Sacrificing participation for safety
- III. Shared Stagnation: Sacrificing progress for participation

[The Solution: Collective Intelligence R&D](#)

- I. The CI Stack: Building the institutions of the future
- II. Value elicitation: Surfacing, aggregating, and understanding conflicting values
- III. Remaking technology institutions: Executing on values via aligned institutions

[Towards a Collectively-Intelligent Future](#)

[Join Us](#)

PDF Version 

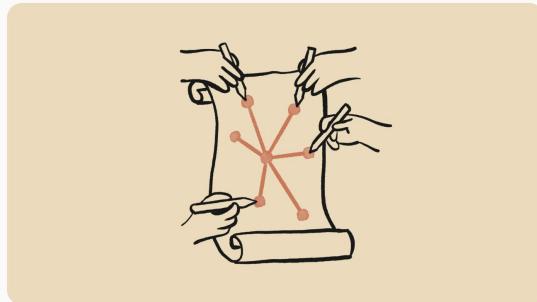
ANTHROPIC

Product Research Company News Careers

Research Societal Impact Policy

Collective Constitutional AI: Aligning a Language Model with Public Input

17 Oct 2023

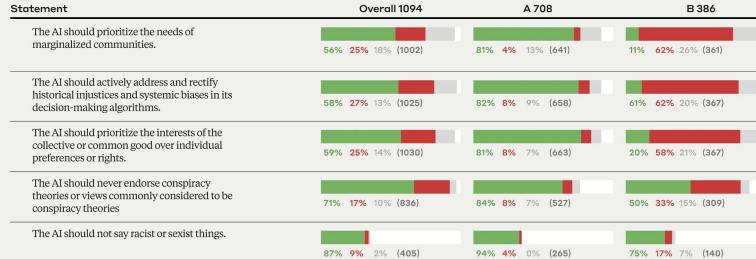


Collective Constitutional AI (CCAI)

A Example of Polis public input process

Group A: 708 participants

Statements which make this group unique, by their votes:



Group B: 386 participants

Statements which make this group unique, by their votes:



<https://www.anthropic.com/news/collective-constitutional-ai-aligning-a-language-model-with-public-input>

Help us pick rules for our AI chatbot!

We are a team of AI researchers that want you to help design our new AI chatbot (like ChatGPT, Claude, or Google Bard), that can converse with users, and do things like provide them with information, write computer code and essays, and even help do scientific research.

Help us pick rules/behavior for our AI. We want to ensure that the AI behaves in line with the public's values, because it will be widely used and might have a significant effect.

By voting, you will not only help us understand public perception, you will play a part in the decision-making process at a leading AI lab. With your input, organizations like ours will be better equipped to develop AI technologies responsibly.

How to participate:

Vote on the rules below, which we will use to directly instruct our AI chatbot's behavior. These are contributed by people like you. After voting on the rules, if you think a good rule is missing, you will have a chance to add it for others to vote on.

You can finish the survey after you have voted on 40 rules. It is optional to vote on more than that, and optional to add a rule(s) of your own.

What rules should our AI follow?

Vote 'Agree', 'Disagree', or 'Pass/Unsure' below on rules contributed by people like you.

Anonymous wrote 100% reading
AI should not discriminate on race or sexual preference
 Agree Disagree Pass/Unsure

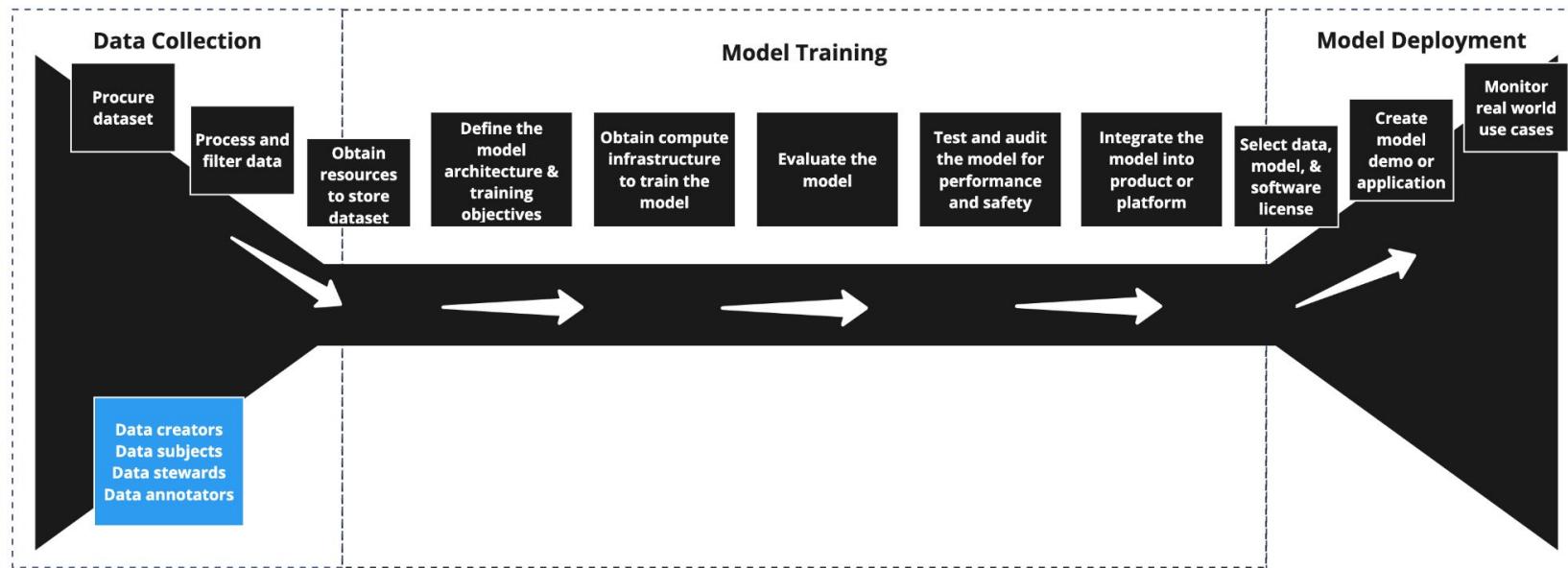
Public Constitution model	Standard Constitution model	Claude Instant 1.2
MMLU Accuracy (%)	72.3	72.4
GSM8K Accuracy (%)	85.6	85.21

Agenda

1. Motivations
2. Case Studies: BigCode, Aya, CIP
3. ***Current Work & Emerging Projects***

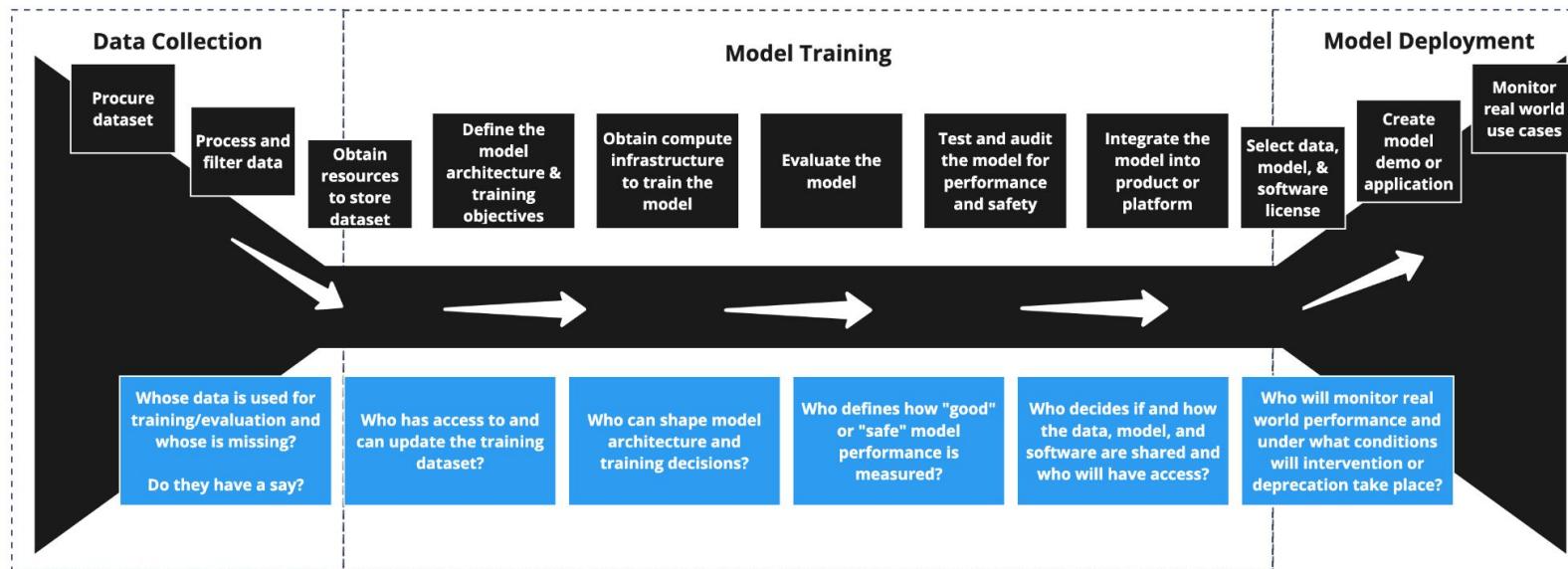
Community-Led LLM Governance

Community governance to empower more people involved in the ML pipeline and ensure more effective, sustainable, and equitable practices and outcomes



Community-Led LLM Governance

Community governance to empower more people involved in the ML pipeline and ensure more effective, sustainable, and equitable practices and outcomes



Emerging Work

- **Collective intelligence to shape data governance resources** for volunteer contributors, to inform future design of volunteer-led data projects and create resources like a data charter template to codify protections and expectations from the volunteer teams.
- **Collective definition of model evaluation priorities**, identifying evaluation priorities and gaps in the existing benchmarking space to shape training and performance expectations from community groups building the training and evaluation datasets

The Alan Turing Institute

Thank you!

Learn more:

- Open Source Initiative 'Defining Open Source AI' Deep Dive:
<https://opensource.org/events/deep-dive-ai-webinar-series-2023>
- Neurips 'Open Source AI Collaborative' Workshop Paper:
<https://arxiv.org/abs/2301.08488>
- BigCode Data Governance Paper: <https://arxiv.org/pdf/2312.03872.pdf>
- Cohere for AI Aya Initiative Data Paper:
<https://cohere.com/research/papers/aya-dataset-paper-2024-02-13>
- CIP whitepaper: <https://cip.org/whitepaper>

Continue the conversation:

- Email jding@turing.ac.uk

