

# Designing for Data Rights in the AI Production Pipeline

Facilitators:

**Jennifer Ding**, The Alan Turing Institute

**Anne Lee Steele**, The Alan Turing Institute

**Yacine Jernite**, Hugging Face

To continue the conversation after today's session, join  
*The Turing Way* Slack or a future community call!

<https://the-turing-way.start.page/>

MozFest 2023 // Allies in Practice

# ⚡ Brain Freeze?



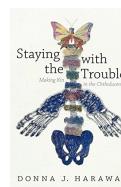
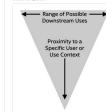
[www.adalovelaceinstitute.org](http://www.adalovelaceinstitute.org)  
Working with the CARE principles:  
operationalising  
Indigenous data  
governance  
Shifting the focus of data governance  
from control to collaboration  
relationships to promote equitable  
and informed participation in data  
processes.



Here are a few resources to get you started

Anything you'd like to share on this topic with the group?

Drop it here!



huggingface.co

Am I in The Stack? - a Hugging Face Space by bigcode

Discover amazing ML apps made by the community



huggingface.co

PII Anonymization - a Hugging Face Space by bigcode

Discover amazing ML apps made by the community



## Case Studies: Choosing an ML License

Below are two hypothetical case studies based on concerns raised through the license development process. While BigScience, a one-year long research workshop on large multilingual models and datasets, Cidney, a hypothetical ML researcher, is working for a...



The BigScience OpenRAIL-M License - Responsible AI Licenses (RAIL)

To help the AI community build the responsible AI (RAI) License more broadly for distributing models, we adopted the Responsible AI License applicable to any associated AI Model.



## deliverypdf.ssrn.com

### The Steep Cost of Capture

6 Pages Posted: 20 Jun 2022 Date Written: 2021 In considering how to tackle this onslaught of industrial AI, we must first recognize that the "advances" in AI celebrated over the past decade were not due to fundamental scientific breakthroughs in AI tec...

# agenda

<b>welcome &amp; intros</b>	00 — 05
<b>activity 1: who &amp; what of ai production</b>	05 — 15
<b>the BigCode approach</b>	15 — 25
<b>activity 2: re-imagining the pipeline</b>	25 — 50
<b>closing</b>	50 — 59

1. Treat yourselves and each other with respect and kindness, per Mozilla's [Community Participation Guidelines](#)
2. Please raise your hand or use the chat for any questions or comments you have during the presentation
3. Have fun! 

# welcome!

Double click on a sticky and share with the group **something about yourself** (e.g. name, pronouns, location)

Jen  
she/her  
London

Nasreen  
She/her  
Ottawa,  
Canada

Remmelt  
He/Him  
Holland

Rich  
He/Him  
Medford,  
MA USA



Anne  
she/her  
London, UK

Evelina  
she/her  
Norway

Mark  
(he/him)  
Chicago, IL  
USA



Yacine  
he/him  
Brooklyn,  
USA

Nadia  
She her  
Madrid

Ushnish  
Sengupta  
he/him  
Oakville,  
Canada



▼  
Dirk  
he/him  
Stroud, UK

Stefan  
he/him  
Berlin,  
Germany

Bushra  
she/her  
Canada



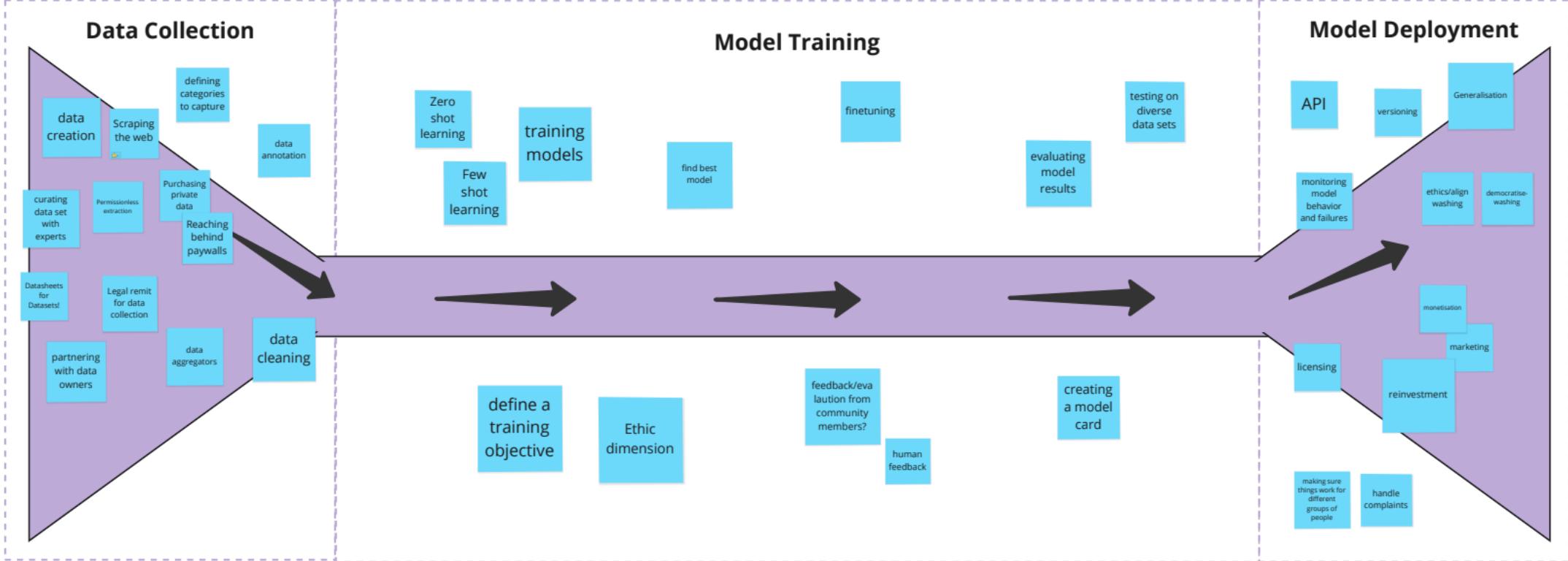
Amelia  
she/they  
Ottawa,  
Canada

Hari Sood  
(he/him),  
London



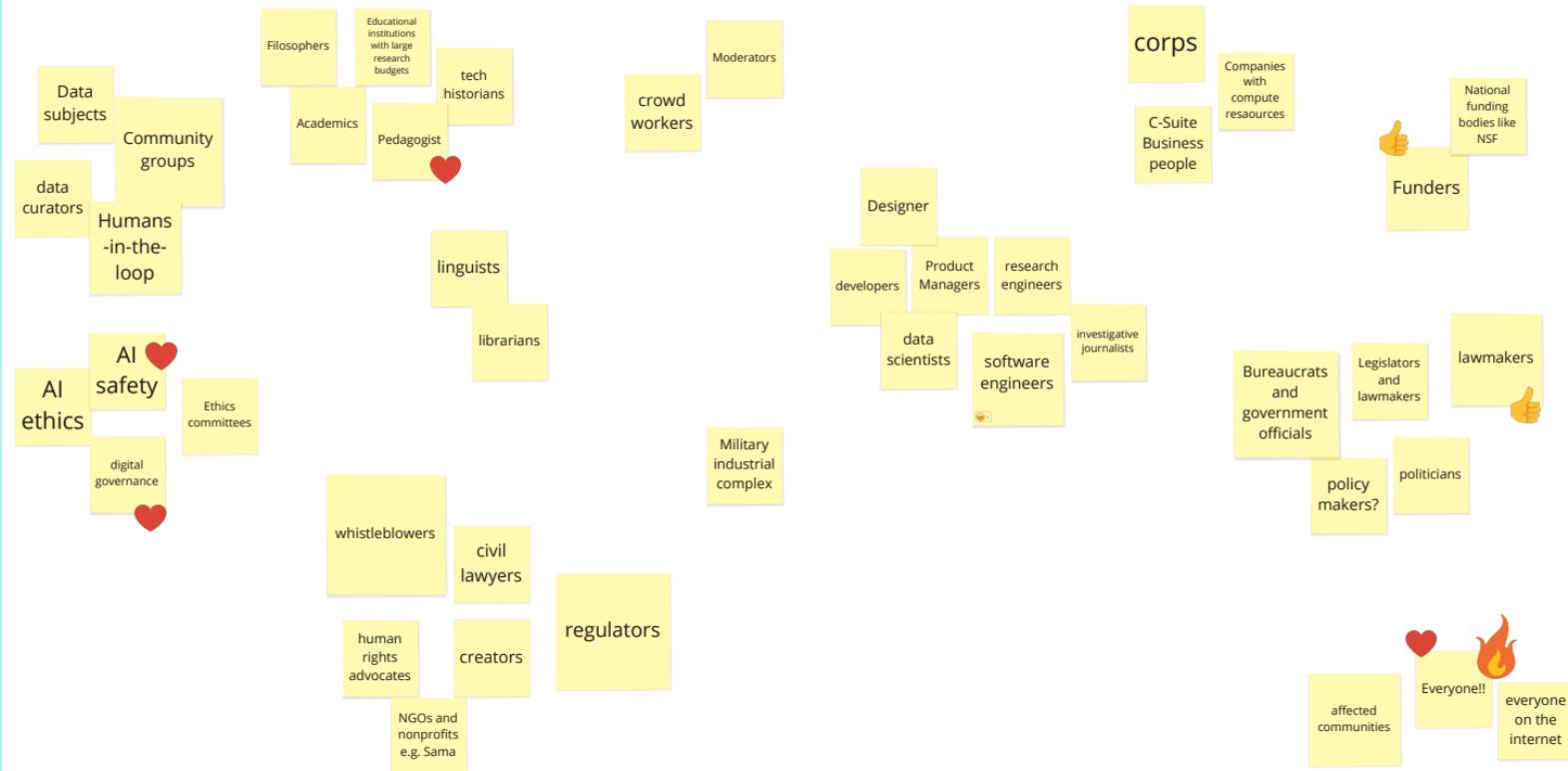
# What kinds of activities are a part of the AI production pipeline?

ai production pipeline



# Who is part of or can shape the AI production pipeline?

who



Drag and duplicate if you like  
what someone else wrote



# Who is impacted by the AI production pipeline?

who

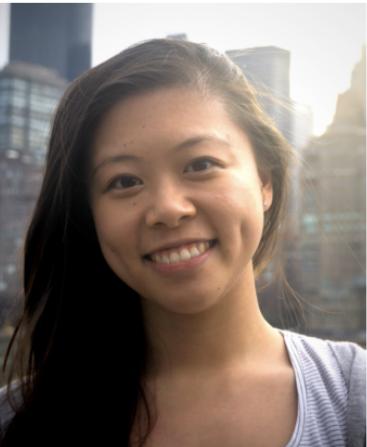


Drag and duplicate if you like  
what someone else wrote



# intros

# Who are we?



Position  
**Research Application Manager**



Position  
**Community Manager, The Turing Way**



Position  
**Research Scientist, Hugging Face**

## Other Resources

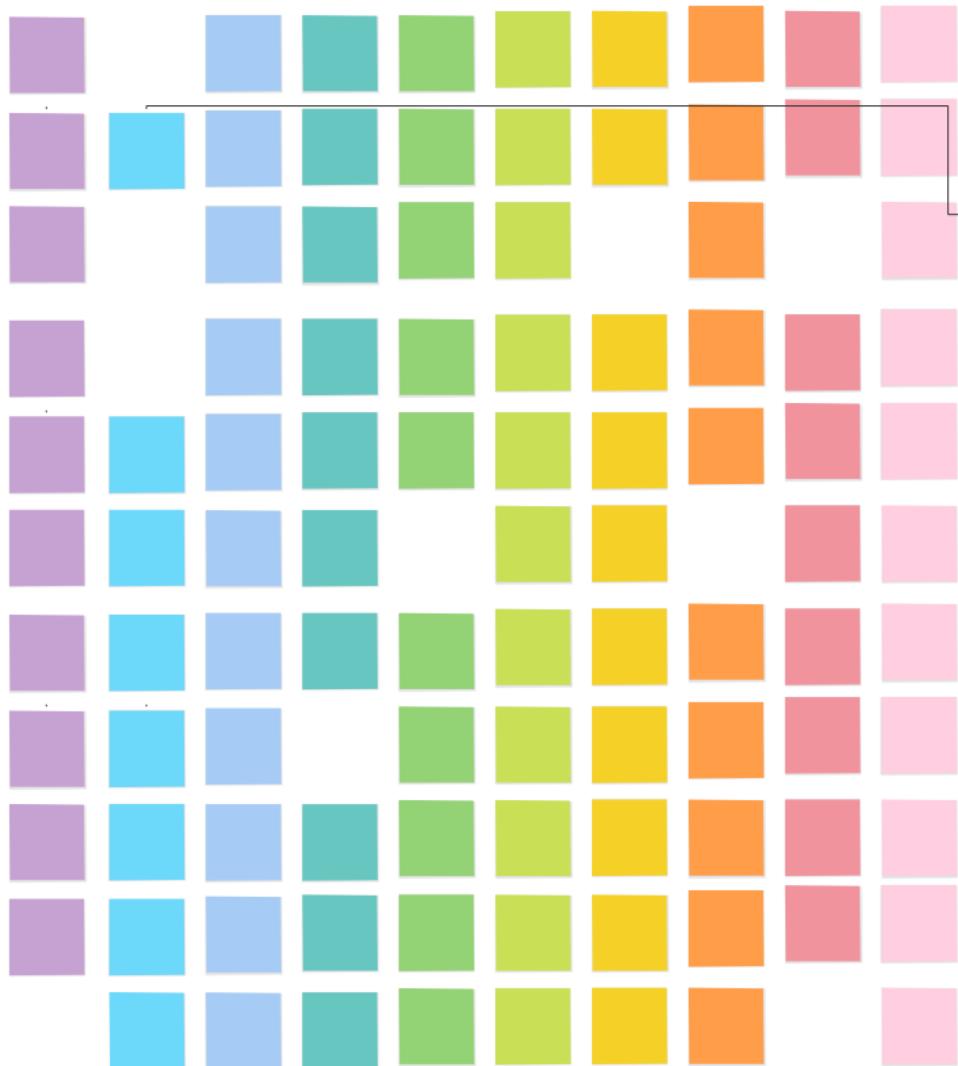
Padlet: <https://pad.sfconservancy.org/p/ttw-mozfest-data-pipeline>

Session info:

<https://schedule.mozillaestival.org/session/KAS9YF-1>

## Sticky-landia

Grab a sticky if you need one!



## case study: BigCode



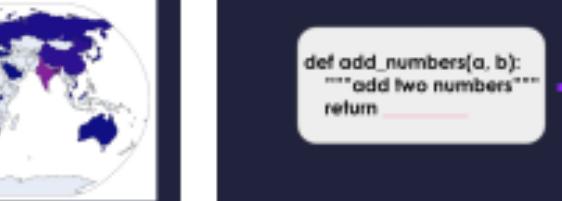
**Hugging Face** **servicenow.**

Open and responsible research & development of large language models for code

### BigCode: Open collaboration

We are building LLMs for code in a collaborative way:

- 500+ participants
- 30+ countries



### What is a language model?

My \_\_\_ → name  
My name \_\_\_ → is  
My name is \_\_\_ → Loubna  
My name is Loubna \_\_\_ → and  
My name is Loubna and \_\_\_ → I  
My name is Loubna and I \_\_\_ → work  
My name is Loubna and I work \_\_\_ → at  
My name is Loubna and I work at \_\_\_ → Hugging Face

### Language models for code generation

Sequence to sequence, encoders, and decoders.



### The Stack

A dataset with **6.4TB** of permissively licensed code in **358 programming languages** with an **Opt-Out** mechanism



<https://huggingface.co/spaces/bigcodein-the-stack>

### Other Data Governance efforts:

#### PII redaction

- Removal of Emails, Keys and IP addresses with regexes
- Next pipeline: PII dataset annotation + train NER model on 7 PII entities for code.

#### Model license

- CodeML OpenRAIL-M v0.1

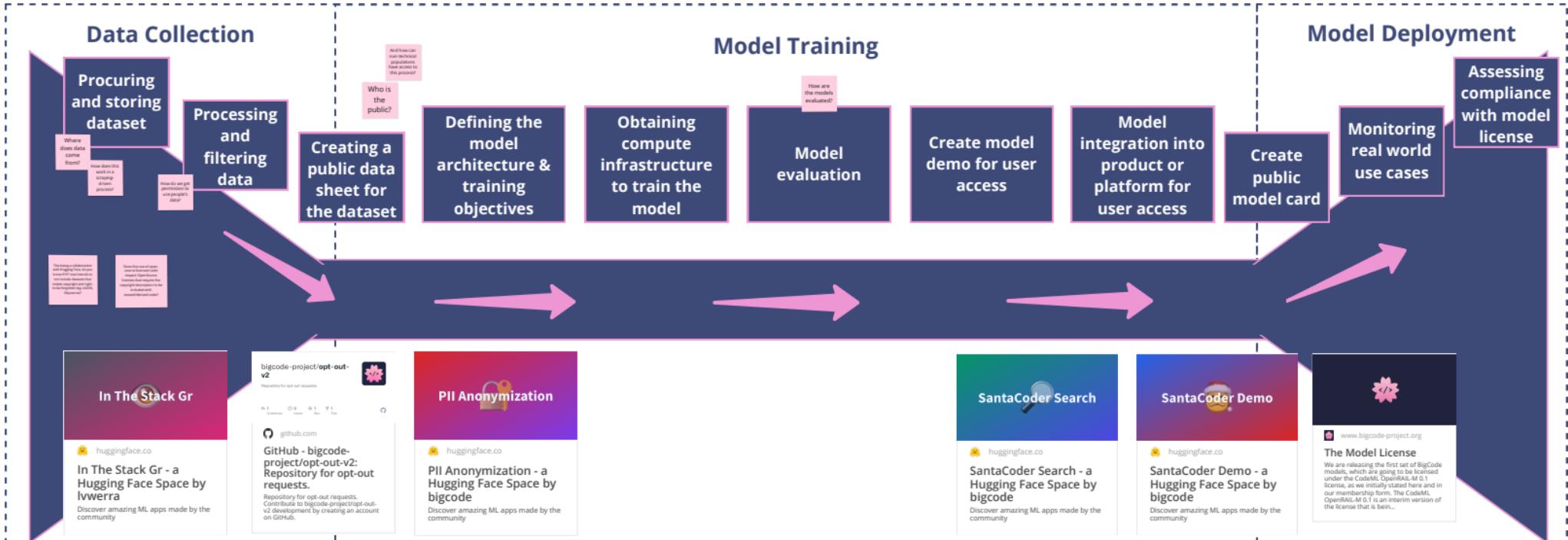
**ai production:** process of building an ML model

**ai pipeline:** steps that form ai production, along which data moves and takes different forms

**data rights:** rights afforded to an individual or community relating to data they create or data that is about them

# BigCode Pipeline

case study: BigCode



What questions, comments, or concerns do you have about the BigCode Pipeline?

# What kinds of activities *should* be a part of the AI production pipeline for/with

Future generations  
?

Re-imagining the AI Production Pipeline

## Data Collection

Be careful with collecting images of children over time. (i.e., predicting how people will age)

ability for data on them to "die" or be removed

History from many point of views

## Model Training

report on carbon emissions / energy costs of training models

research: fundamental limits of controllability of system effects (measurability, modellability, predictability, error-correctability)

## Model Deployment

use compute powered by renewable energy

algorithmic disgorgement: destroy models companies that scale up large spaghetti code for general use

# What kinds of activities *should* be a part of the AI production pipeline for/with

## Re-imagining the AI Production Pipeline

### Data Collection

stricter data collection laws and regulations

community-based data collection agreements

Avoid further reinforcing stereotypes

### Model Training

models being developed in partnership with communities

Make sure the model isn't overly biased.

ability to interrogate models easily / audit how they work

### Model Deployment

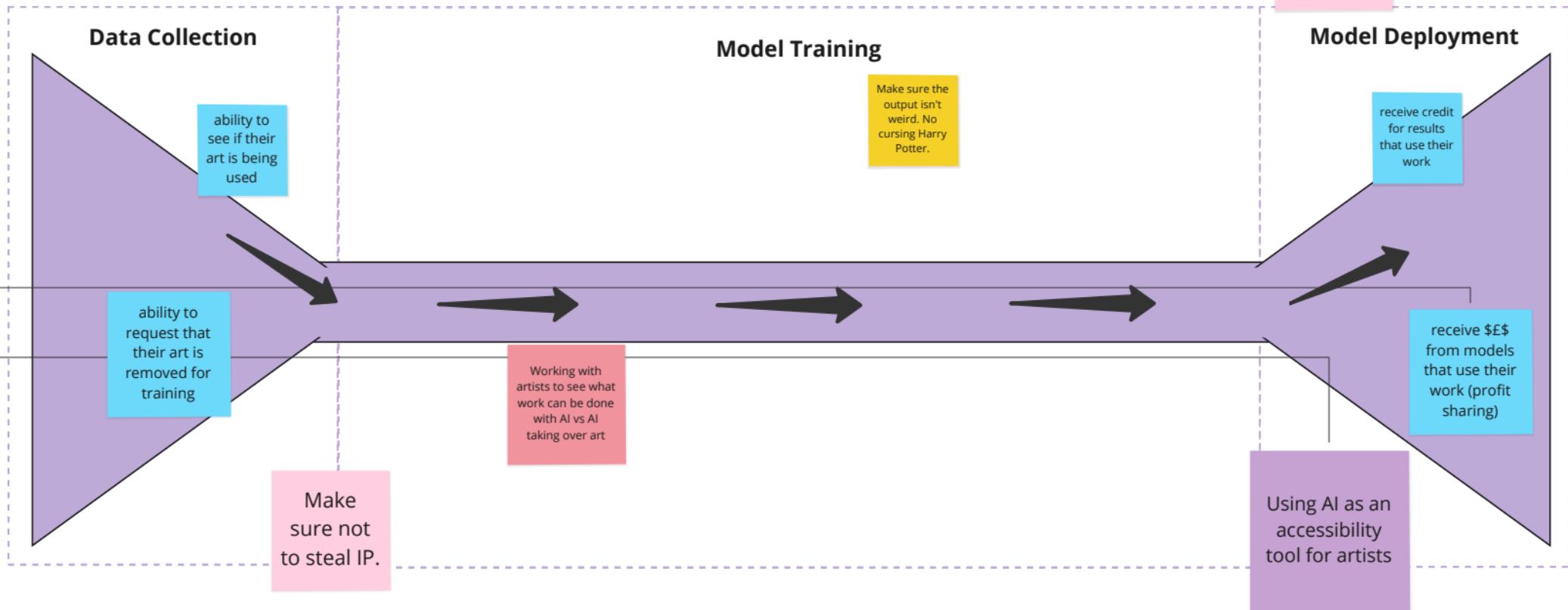
require poorly performing models to be removed from operations

Policed populations

?

# What kinds of activities *should* be a part of the AI production pipeline for/with Artists ?

## Re-imagining the AI Production Pipeline



# What kinds of activities *should* be a part of the AI production pipeline for/with

People in  
global south  
who live  
outside net of  
AI

?

Re-imagining the AI Production Pipeline

## Data Collection

right to  
shape what  
data is /  
isn't used

History from  
many point of  
views not only the  
white European  
colonizer

## Model Training

inform  
data  
selection

Be part of the  
production and  
design, not just  
data subjects  
and annotators

## Model Deployment

Make sure  
it benefits  
the  
population

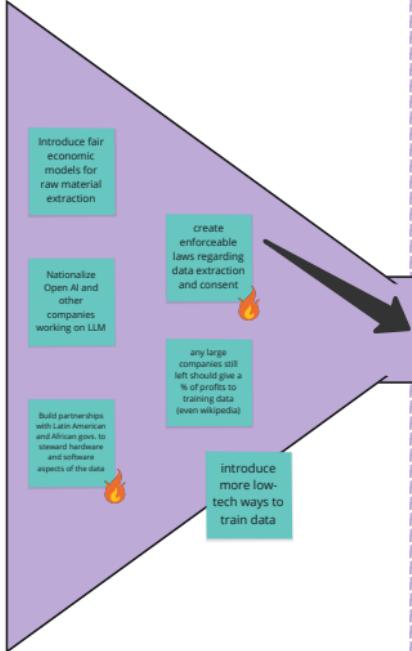
guaranteed  
access for  
the data  
subjects

# What kinds of activities *should* be a part of the AI production pipeline for/with

miners ?

## Re-imagining the AI Production Pipeline

### Data Collection



### Model Training

limit compute access to public utilities

allow miners to be a part of training process

### Model Deployment

Miners, and other global south ppl become PMs and developers

Have ability to decide use-cases, features, etc.

AI towards utility of the masses as opposed to a utility of the silicon valley class

What kinds of activities *should* be a part of the AI production pipeline for/with

XXX

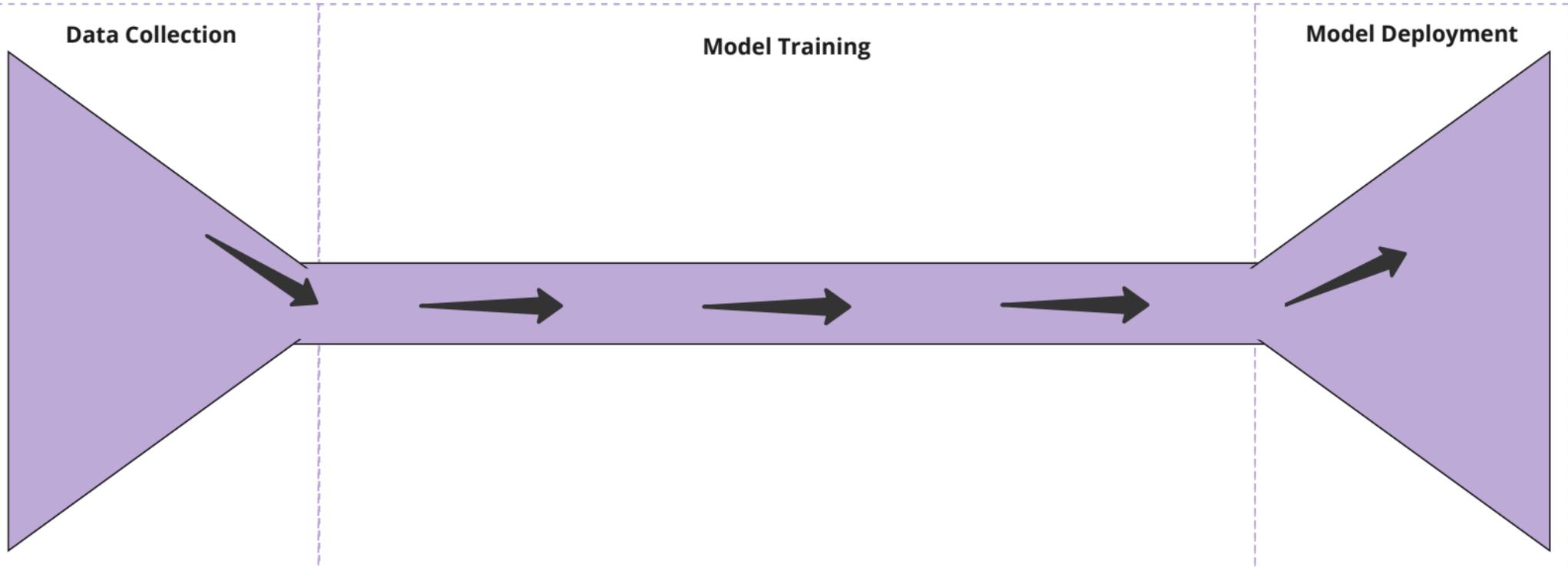
?

Re-imagining the AI Production Pipeline

Data Collection

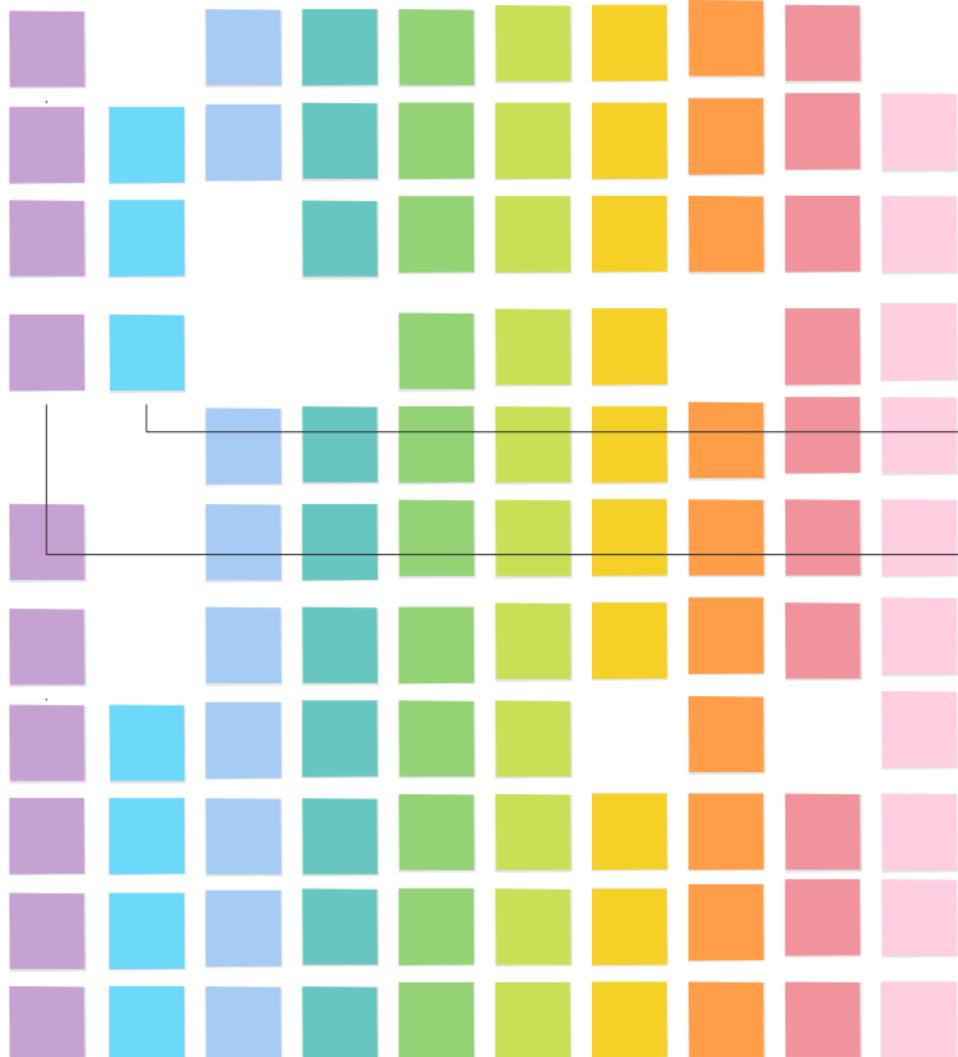
Model Training

Model Deployment



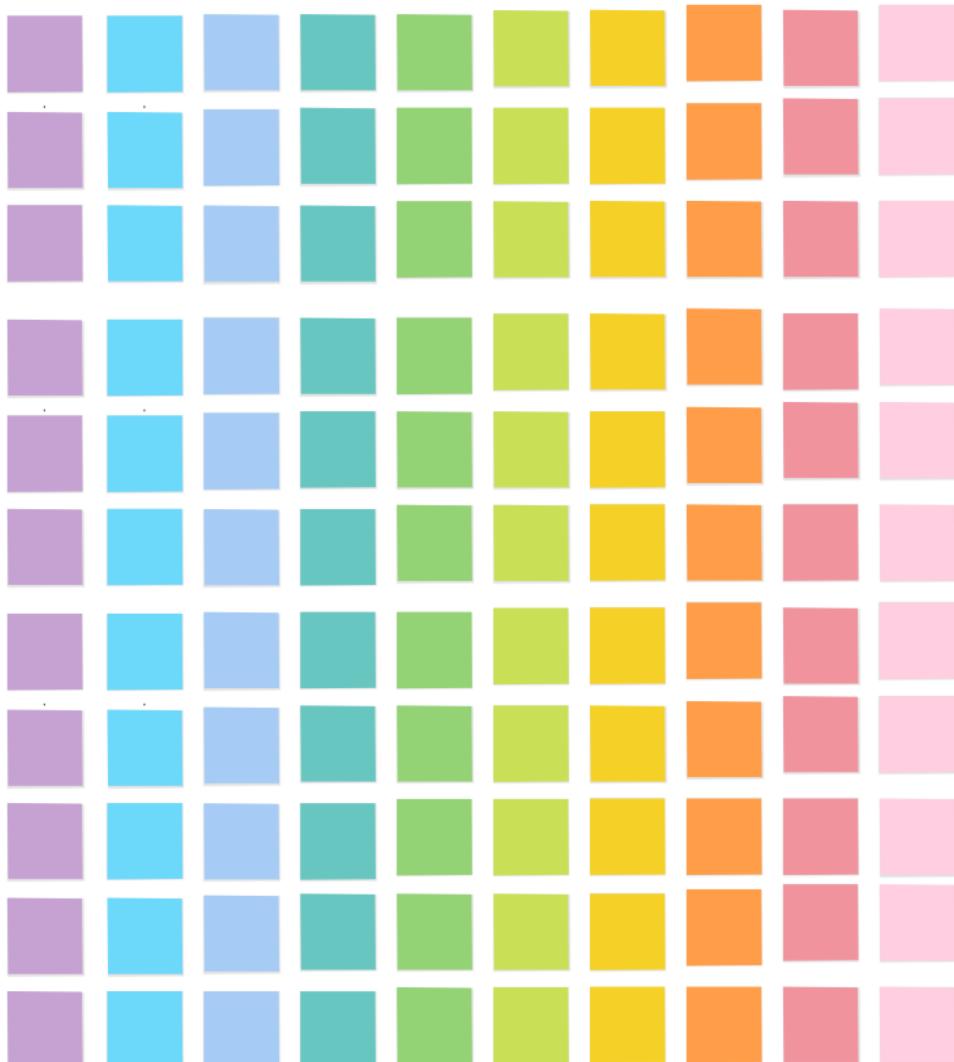
## Sticky-landia

Grab a sticky if you need one!



## Sticky-landia

Grab a sticky if you need one!



# final announcements

## bigcode

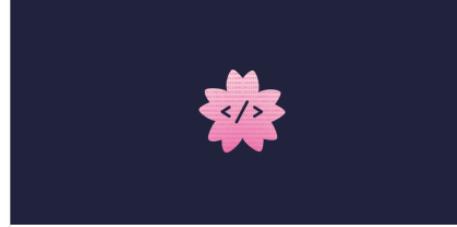


A screenshot of a GitHub repository page for "bigcode". The repository icon is a purple gear with white symbols. The repository name "bigcode" is displayed below it. At the bottom of the page, there is a link "huggingface.co/bigcode" with a yellow emoji icon.

 [huggingface.co](https://huggingface.co/bigcode)

### bigcode (BigCode)

We're on a journey to advance and democratize artificial intelligence through open source and open science.



A screenshot of the BigCode Project website. The header features a large purple flower icon with a pink center containing the code symbol "</>". Below the header, there is a link "www.bigcode-project.org" with a purple emoji icon.

### Open and responsible development of LLMs for code

BigCode is an open scientific collaboration working on the responsible development of large language models for code

# final announcements

## the turing way

**Reflect, Unlearn, Reframe:  
Community Care in times of  
digital burnout**

31 March, 14:00 - 15:30 UTC+1 – [Register on Eventbrite](#)

Agnes Kirraga  
Technical Lead  
African Population Health Research Council

Chris Hartgerink  
CEO & Founder  
Liberate Science

Maya Sundukova  
Resident Fellow  
Open Life Science

Patricia Herterich  
Resident Fellow  
Open Life Science

Eirini Zormpa  
Community Manager  
AIM-RSF & The Turing Way

Hosted by: Open Life Science

[www.eventbrite.co.uk](#)

Reflect, unlearn, reframe: Community care in times of digital burnout

The Turing Way Fireside Chat series features people, ideas and projects in open and reproducible research.

Click here to see links and & resources about the project

[the-turing-way.start.page](#)

### Welcome to The Turing Way!

Check out this welcome page for more information about our guides, how to contribute, community calls, upcoming events, and more!

...Or at a Collaboration Cafe!



[hacmd.io](#)  
.The Turing Way,  
online Collaboration  
Cafe | 5 April 2023 -  
HackMD

Join us at our  
upcoming  
Fireside Chat!

MozFest 2023 // Allies in Practice

# Designing for Data Rights in the AI Production Pipeline

Facilitators:

**Jennifer Ding**, The Alan Turing Institute

jding@turing.ac.uk

**Anne Lee Steele**, The Alan Turing Institute

asteele@turing.ac.uk

**Yacine Jernite**, Hugging Face

yacine@huggingface.co

To continue the conversation after today's session, join  
*The Turing Way* Slack or a future community call!

<https://the-turing-way.start.page/>

MozFest 2023 // Allies in Practice

## Data Collection

Procuring and storing dataset

Processing and filtering data

Creating a public data sheet for the dataset

Defining the model architecture & training objectives

Obtaining compute infrastructure to train the model

## Model Training

Model evaluation

Create model demo for user access

Model integration into product or platform for user access

Create public model card

## Model Deployment

Assessing compliance with model license

### In The Stack Gr

In The Stack Gr - a Hugging Face Space by lvverra  
Discover amazing ML apps made by the community



GitHub - bigcode-project/opt-out-v2: Repository for opt-out requests.  
Repository for opt-out requests. Contribute to bigcode-project/opt-out-v2 development by creating an account on GitHub.

### PII Anonymization

PII Anonymization - a Hugging Face Space by bigcode  
Discover amazing ML apps made by the community



### SantaCoder Search

SantaCoder Search - a Hugging Face Space by bigcode  
Discover amazing ML apps made by the community



### SantaCoder Demo

SantaCoder Demo - a Hugging Face Space by bigcode  
Discover amazing ML apps made by the community



### The Model License

We are releasing the first set of BigCode models, which are going to be licensed under the CodeML OpenRAIL-M 0.1 license, as we have done with our membership form. The CodeML OpenRAIL-M 0.1 is an interim version of the license that is bei...

