

ITU, WHO & WIPO Global Initiative on AI for Health (GI-AI4H): Purpose, governance, working & topic groups

Focused deep-dive:

*Data standards for health AI:
Benchmarking, metadata and
federated data discovery*

Overview

ITU/WHO/WIPO

Global Initiative on AI for Health
(GI-AI4H)

Simão Campos

ITU

AI for Good Webinar | 27 March 2026



A large orange circle is positioned on the left side of the slide, partially overlapping the white background.

What is the GI-AI4H?

Joint UN initiative led by ITU, WHO and WIPO to harness AI for the benefit of health.

Global coordination platform for standards, governance, knowledge sharing and implementation support.

Builds on the legacy and deliverables of the ITU/WHO Focus Group on AI for Health (FG-AI4H).

Strategic pillars

Enable

- Development of international standards, norms, guidance and evidence-based technical products.

Facilitate

- Convening stakeholders and operating enabling mechanisms (e.g. Open Code Infrastructure).

Implement

- Supporting country-level scaling, capacity building and sustainable implementation models.

Governance structure

Steering Committee

Senior representatives from ITU, WHO and WIPO.

Joint Secretariat

Operational coordination across the three organizations.

Working Groups, Topic Groups & Facilitation Groups

Partners and donors

Supporting activities and implementation.

How the work is organized

Working Groups (WGs)

Cross-cutting themes such as ethics, regulation, data, evaluation and intellectual property.

Topic Groups (TGs)

Specific health domains (e.g. traditional medicine, maternal health).

Facilitation groups

Stakeholder networks supporting knowledge exchange and implementation.

Current Working & Topic Groups (snapshot)

Working Groups

- WG – Regulatory considerations
- WG – Ethics and governance
- WG – Data
- WG – Evaluation
- WG – Intellectual property & innovation

Topic Groups (tentative)

- Traditional medicine ✓
- Maternal & reproductive health
- Point-of-care / primary health care
- Oral health (currently on hold)

WG-DATA: scope and rationale

Addresses fragmentation and high transaction costs in health-AI data usage.

Develops a Data and Model Exchange Protocol (DMXP) to support discoverability and access.

Builds on FG-AI4H benchmarking experience and modern data standards (e.g., Croissant).

WG-DATA objectives

G1: Dataset metadata standardization and basic matchmaking.

G2: Searchable, ontology-linked index of health datasets.

G3: Secure transaction protocols
(authentication, licensing, access control).

G4: Federated testing environments and real-world pilots.

Data Standards for Health AI

Benchmarking, Metadata and
Federated Data Discovery

Marc Lecoultre

Open Code Infrastructure | GI-AI4H
ITU-WHO-WIPO Global Initiative on AI for Health

AI for Good Webinar | 27 March 2026



The Data Challenge in Health AI

- Health data is siloed across countries, institutions, and systems
- No consistent way to describe, discover, or compare datasets
- Every dataset has its own format, documentation, and access model
- Regulatory requirements vary by jurisdiction

Impact on Health AI

- AI developers cannot find suitable training and benchmarking data
- Regulators cannot verify what data was used to build a model
- Research is duplicated — datasets are not discoverable
- Benchmarking is inconsistent — models tested on incomparable data
- Health equity gaps widen when underrepresented data stays invisible

We need a standardized, federated approach to health dataset discovery

The Open Code Infrastructure: A Public Good

The **AI4H Assessment Platform** is a free, open-access platform built as a public good, providing a secure, end-to-end framework for developing and rigorously assessing health AI algorithms globally.

Open Source

Transparent, auditable,
community-driven

Universal

Across borders,
stakeholders, disciplines

Health-Specific

Privacy, ethics, consent,
regulatory alignment

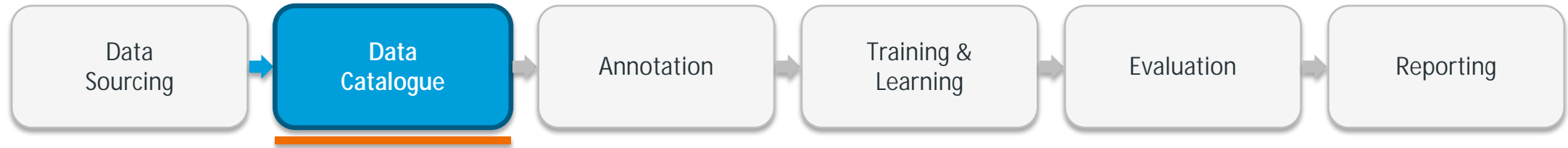
- 10 module drivers, 40+ contributors from 5 continents
- Engineers, academics, medical doctors, regulators
- Integrates data privacy, ethics, governance, and secure metadata management

Data Catalogue

Discovering Health Datasets
for AI Benchmarking



End-to-End Platform: Focus on Data Catalogue



The Data Catalogue is the backbone of health AI benchmarking

- Standardized dataset descriptions using Croissant metadata
- Browse and search datasets by modality, disease, population
- Machine-readable metadata enables automated quality checks
- Federated architecture: data stays at source, metadata travels
- Supports regulatory traceability and audit requirements
- Interoperable with Kaggle, Hugging Face, OpenML, Google Dataset Search

Why Data Standards Matter for Health AI

Without Standards

- ✗ Dataset described in a README — "100 chest X-rays, contact PI for access"
- ✗ No consistent vocabulary across institutions or countries
- ✗ Manual effort to understand format, structure, and licensing
- ✗ Cannot compare datasets used in different benchmarks
- ✗ Regulators have no audit trail for data provenance
- ✗ Discovery relies on personal networks and word-of-mouth

With Croissant

- ✓ Machine-readable JSON-LD metadata indexed by search engines
- ✓ Shared vocabulary linked to biomedical ontologies (SNOMED, ICD, LOINC)
- ✓ Self-describing structure: fields, types, splits, and ML semantics
- ✓ Datasets are comparable across challenges and institutions
- ✓ Full provenance: collection method, annotation, consent, governance
- ✓ Global discovery via federated catalogue and standard APIs

Croissant: The Metadata Standard for ML Datasets

Community built and open standard by MLCommons • Built on schema.org • JSON-LD format

1

Metadata Dataset description, license, creators, responsible AI attributes

2

Resources Pointers to data files, URLs, checksums for integrity

3

Structure Records, fields, data types — how raw data maps to tables

4

ML Semantics Train/test/validation splits, intended tasks and use cases

Adopted by: [Google Dataset Search](#) • [Kaggle](#) • [Hugging Face](#) • [OpenML](#) • [CKAN](#) • [TensorFlow Datasets](#)

WG "BioCroissant" - Extending Croissant for Health & Life Sciences

The life sciences extension of Croissant (v1.1) — domain-specific vocabulary for biomedical datasets

Biomedical Ontologies

Links to SNOMED CT, ICD, LOINC, MeSH, UMLS

Clinical Context

Imaging modality, anatomical region, pathology, cohort

Patient Metadata

Demographics, consent status, cohort characteristics

Regulatory Attributes

IRB approval, data governance, privacy level, jurisdiction

Provenance

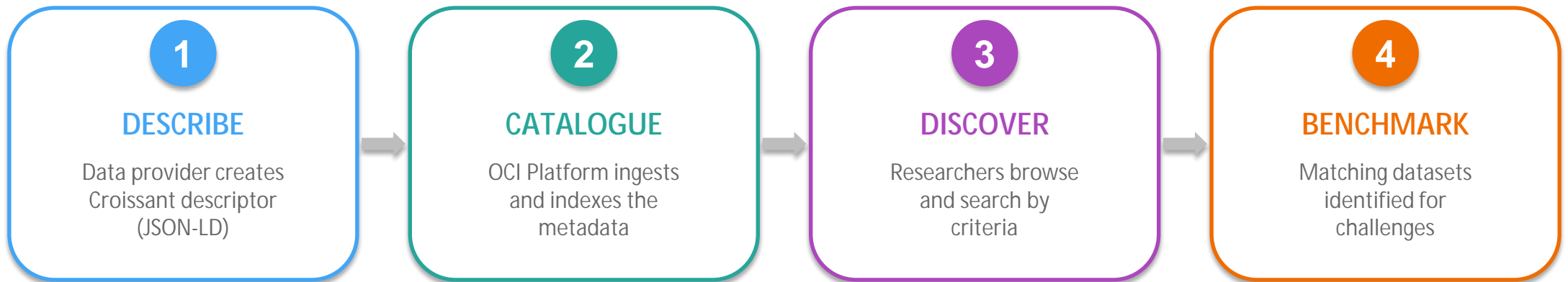
Collection protocol, annotation methodology, ground truth

Cross-border Comparison

Shared vocabulary enables global benchmarking

A chest X-ray dataset from Kenya and one from Germany can be described, compared, and discovered using the same standardized vocabulary

How It Works: Dataset Cataloguing with Croissant



Key Principle: Data stays at source — only metadata travels

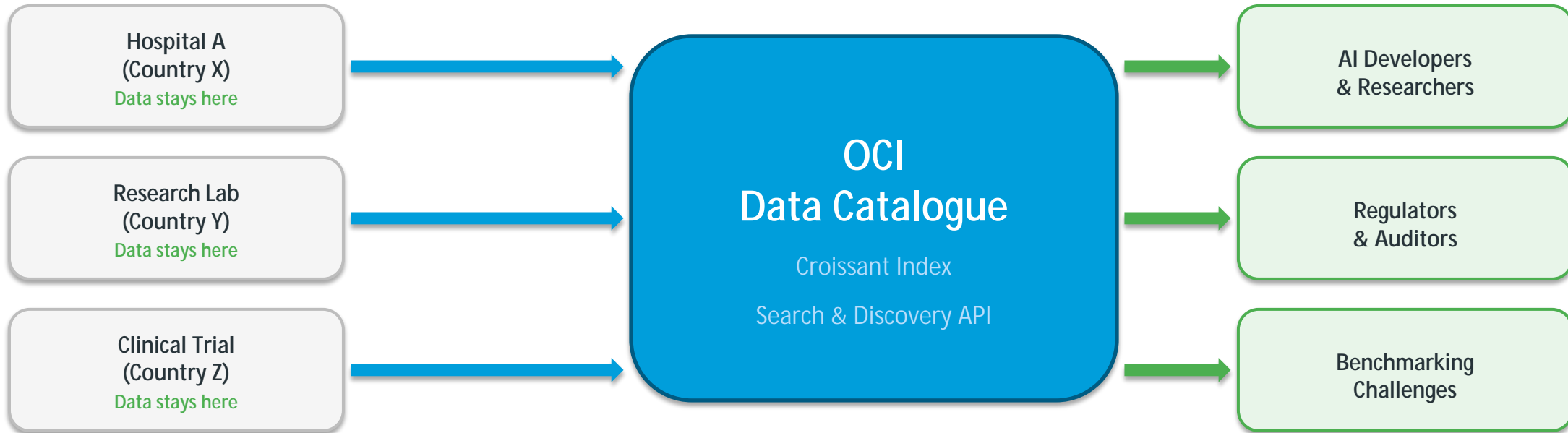
Federated catalogue respects data sovereignty while enabling global discovery

Browse by: disease area • imaging modality • population • geography • consent status • annotation quality

Federated Data Discovery Architecture

Metadata flows →

→ Search & discovery



Data Multiplexers

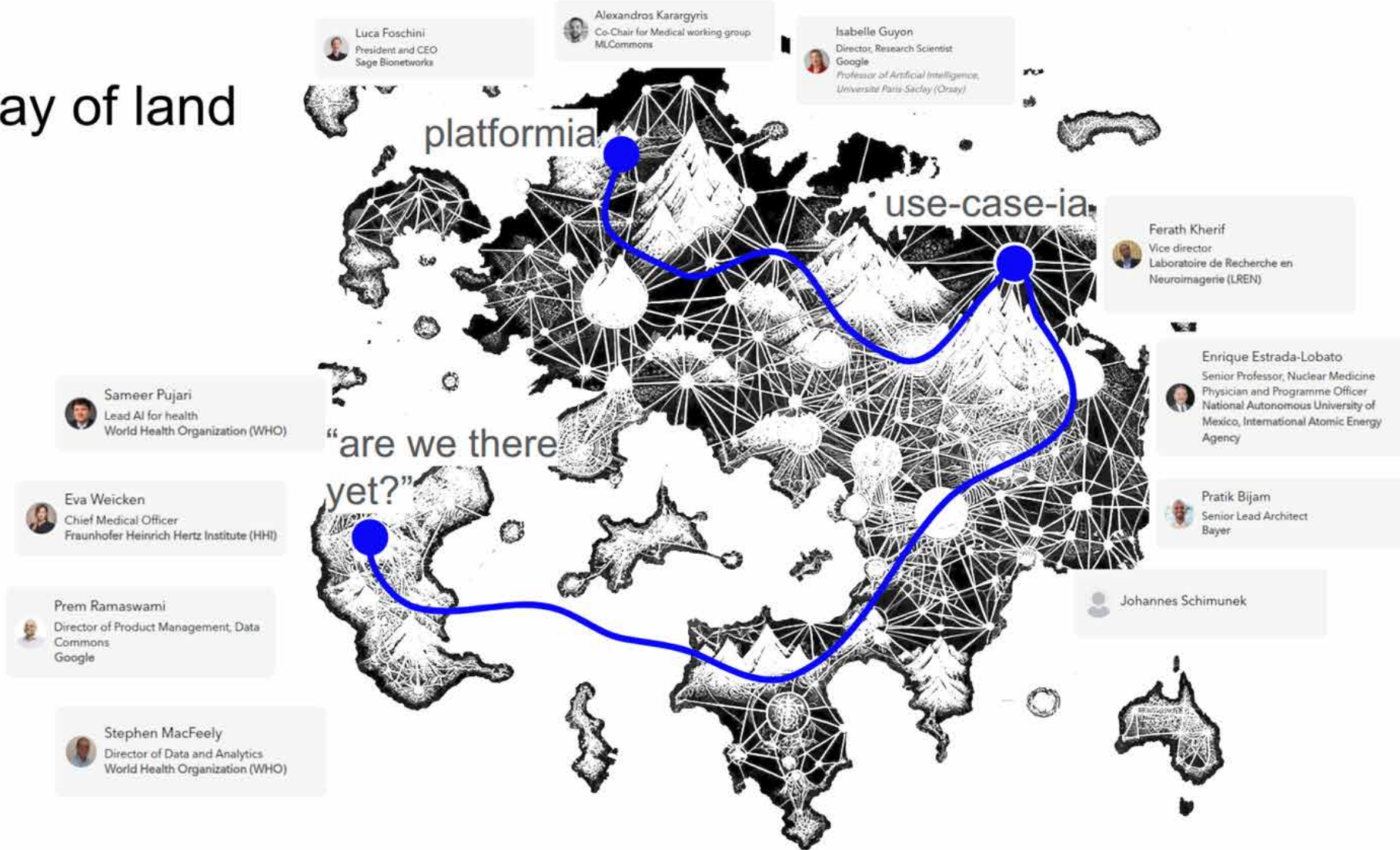
Infrastructure for Open Data Marketplaces
and Exchanges

Luis Oala

ITU-WHO-WIPO Global Initiative on AI for Health
Brickroad

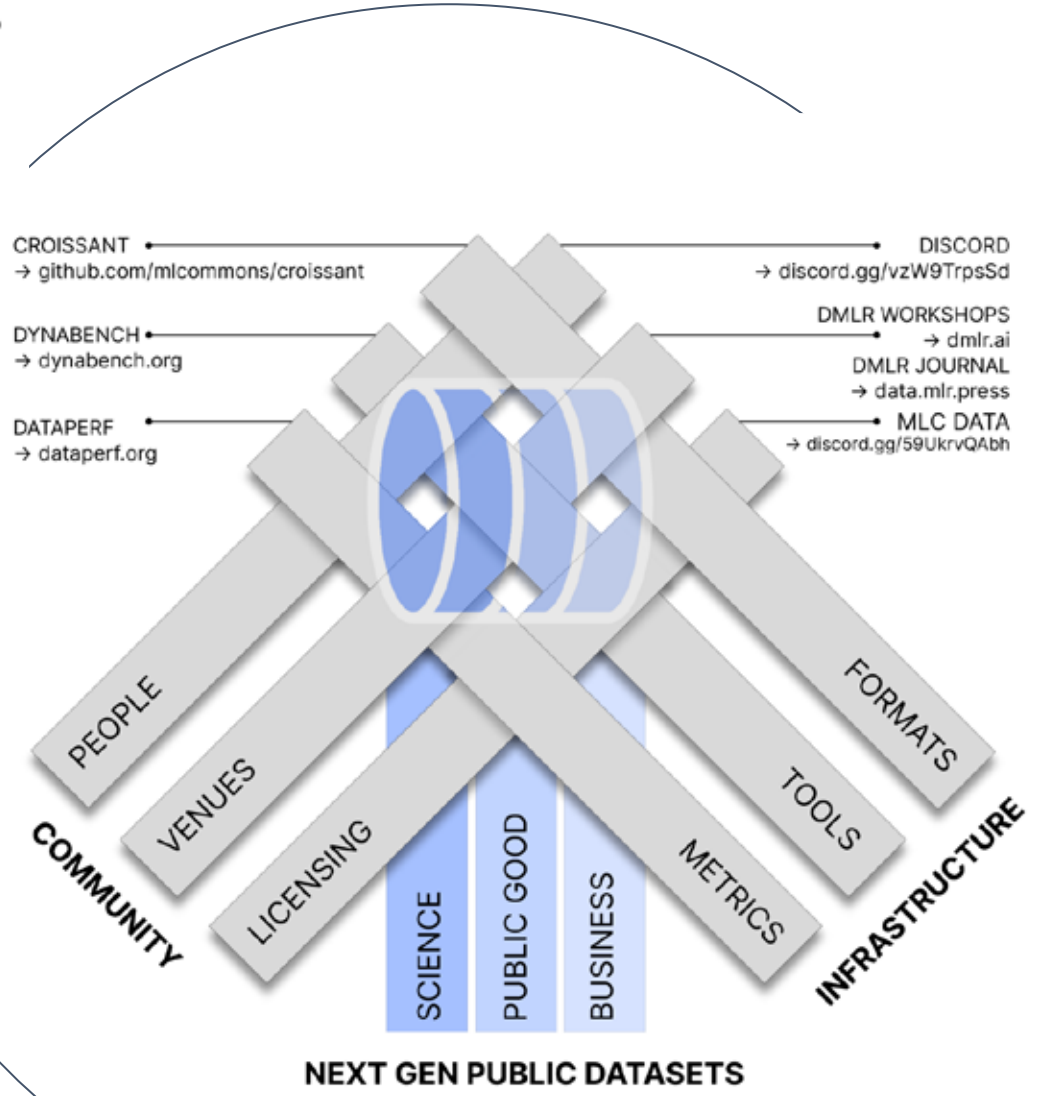
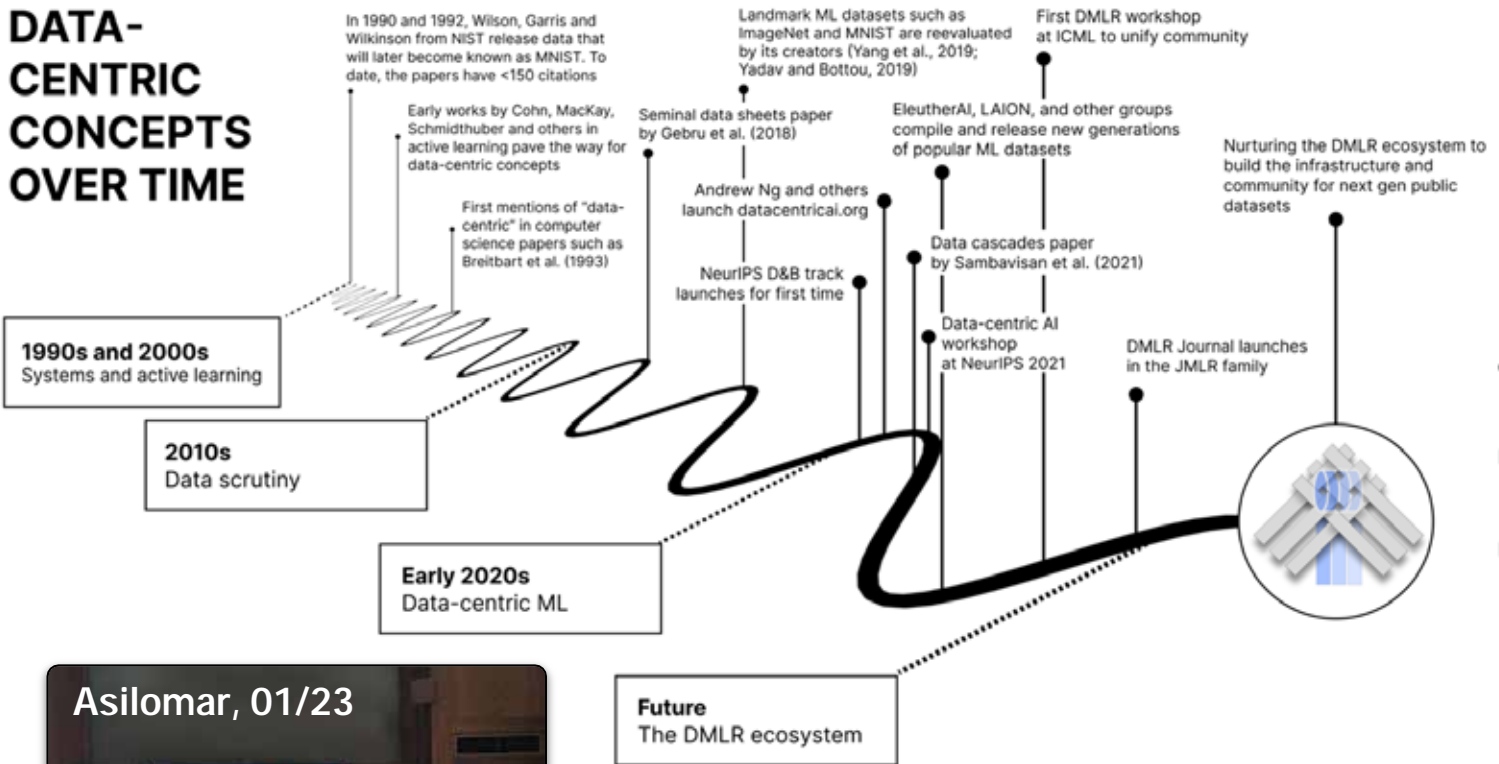


lay of land

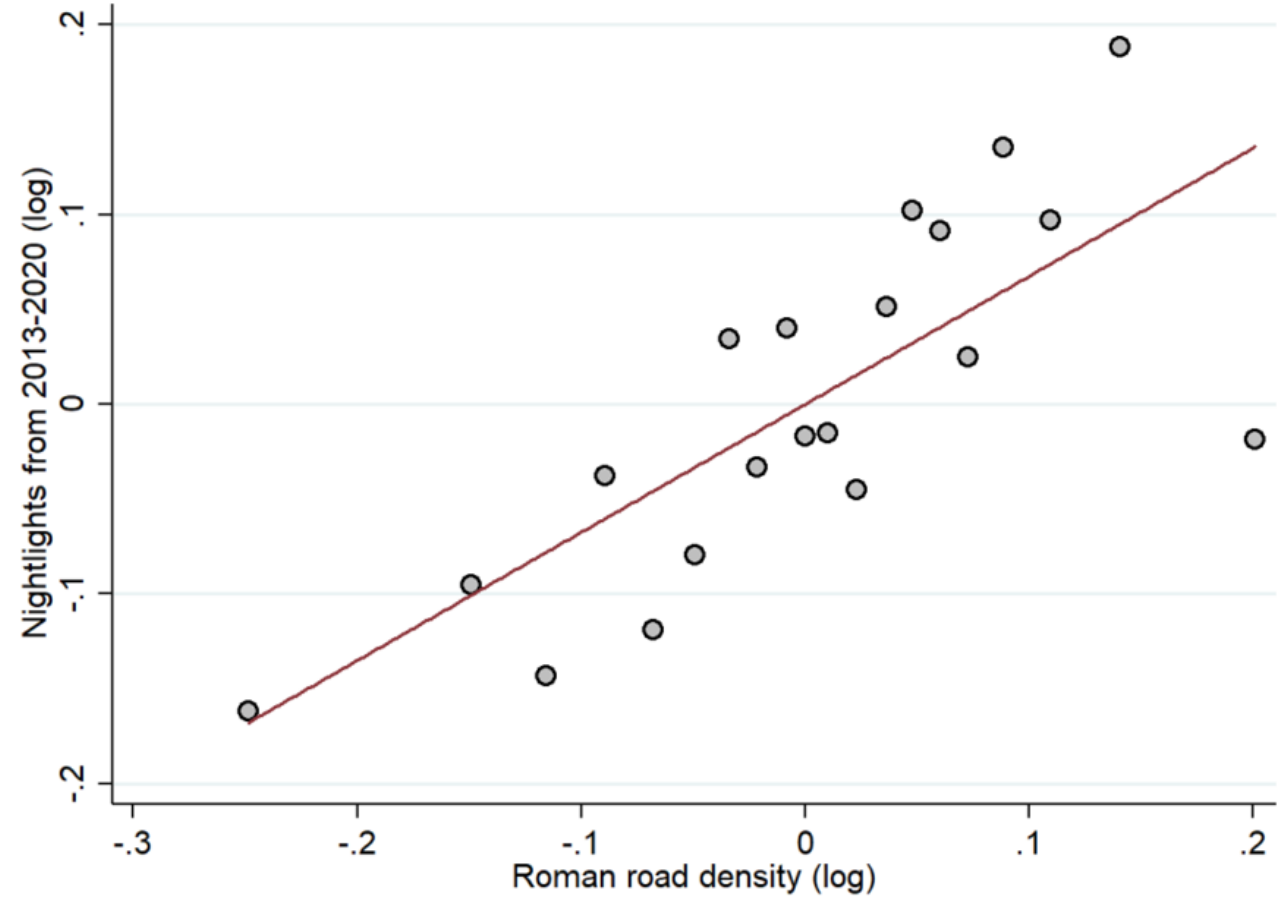


Throwback: Geneva, 31 May 2024

DATA-CENTRIC CONCEPTS OVER TIME

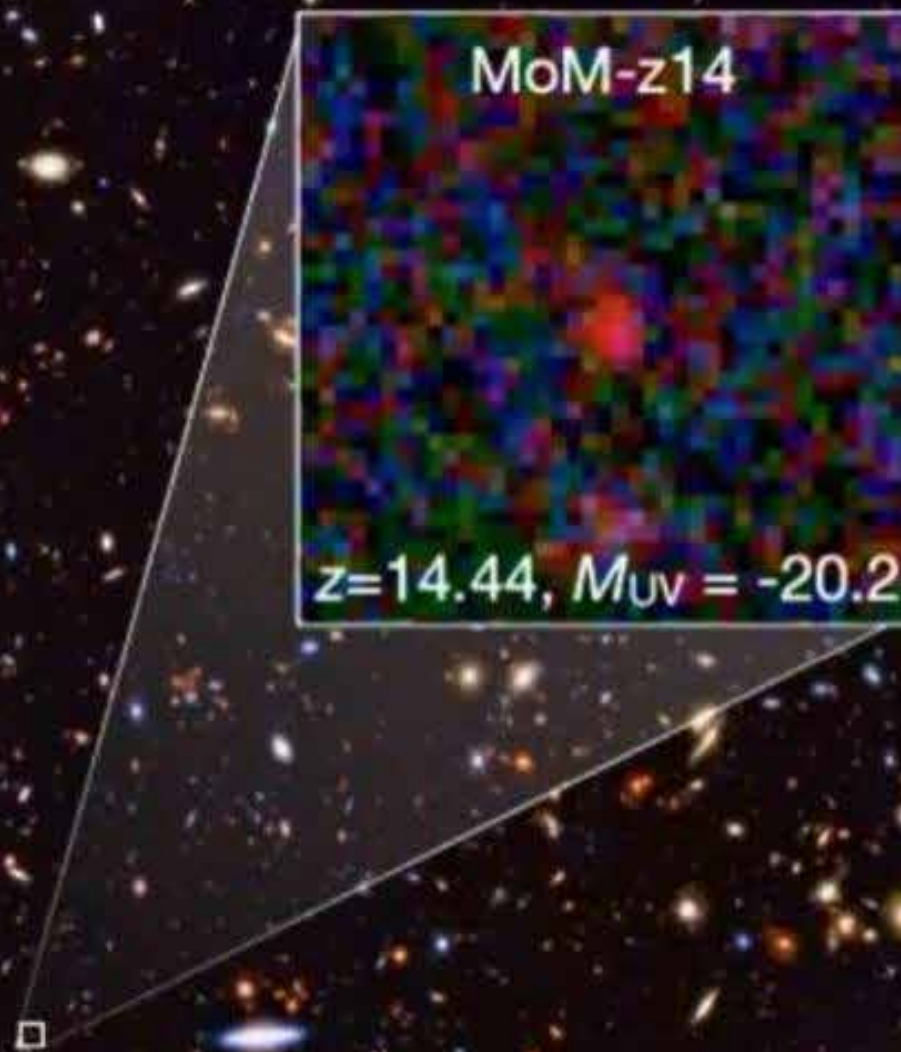


Assumptions



assumptions

- value creation at the information frontier



Assumptions

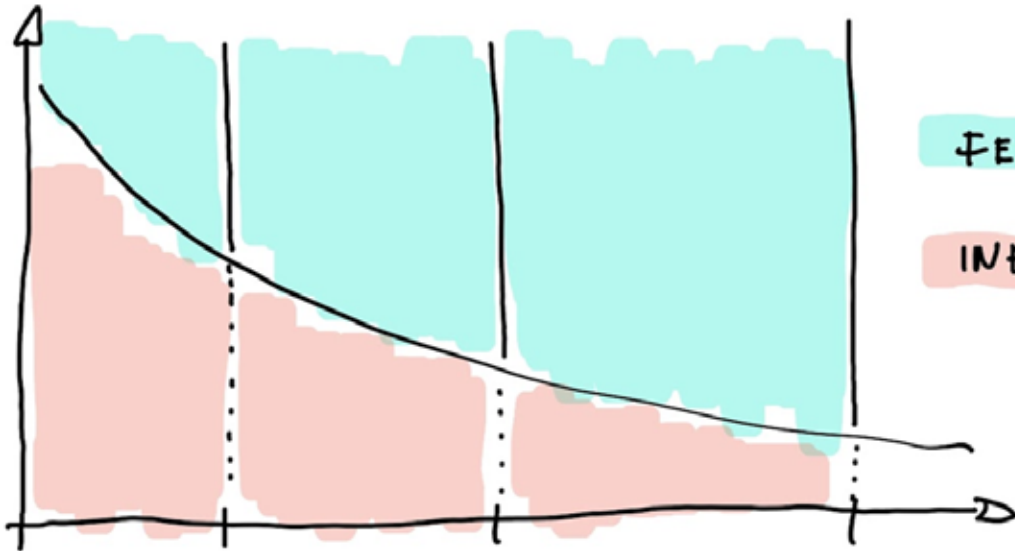
- value creation at the information frontier
- long tail of tasks



Goals

Increase the space of feasible tasks

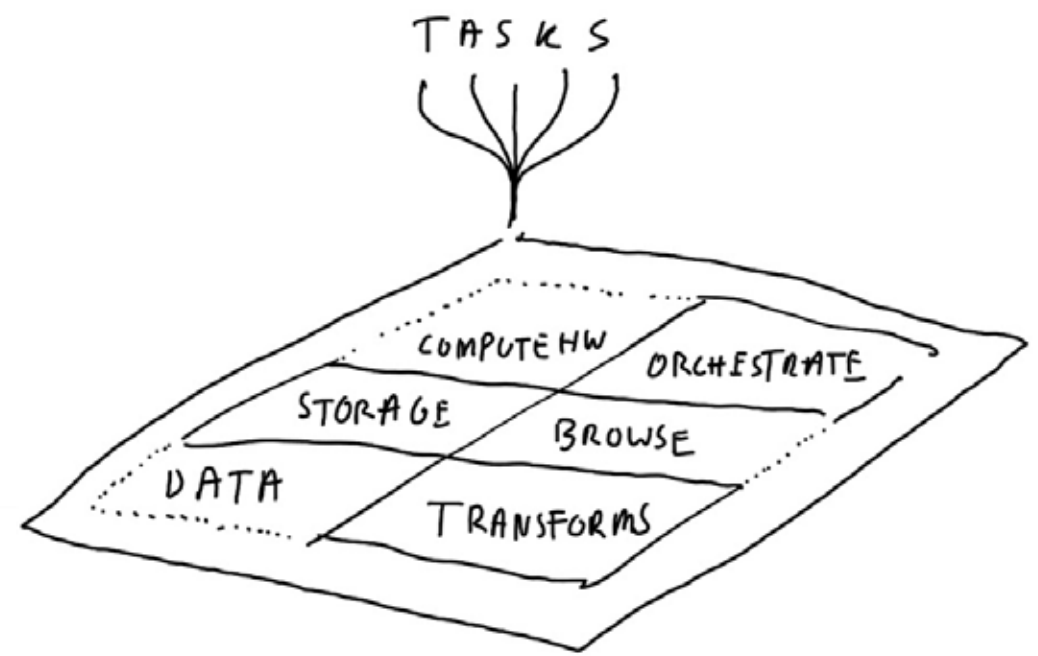
TRANSACTION COSTS IN ML



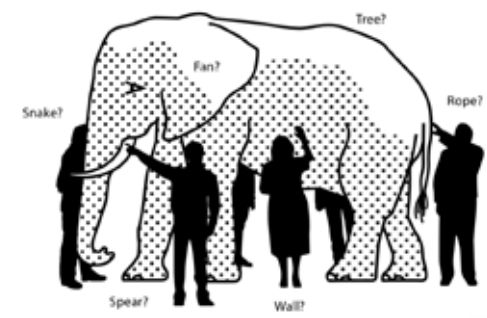
FEASIBLE TASKS

INFEASIBLE TASKS

ML SYSTEMS ADVANCEMENTS



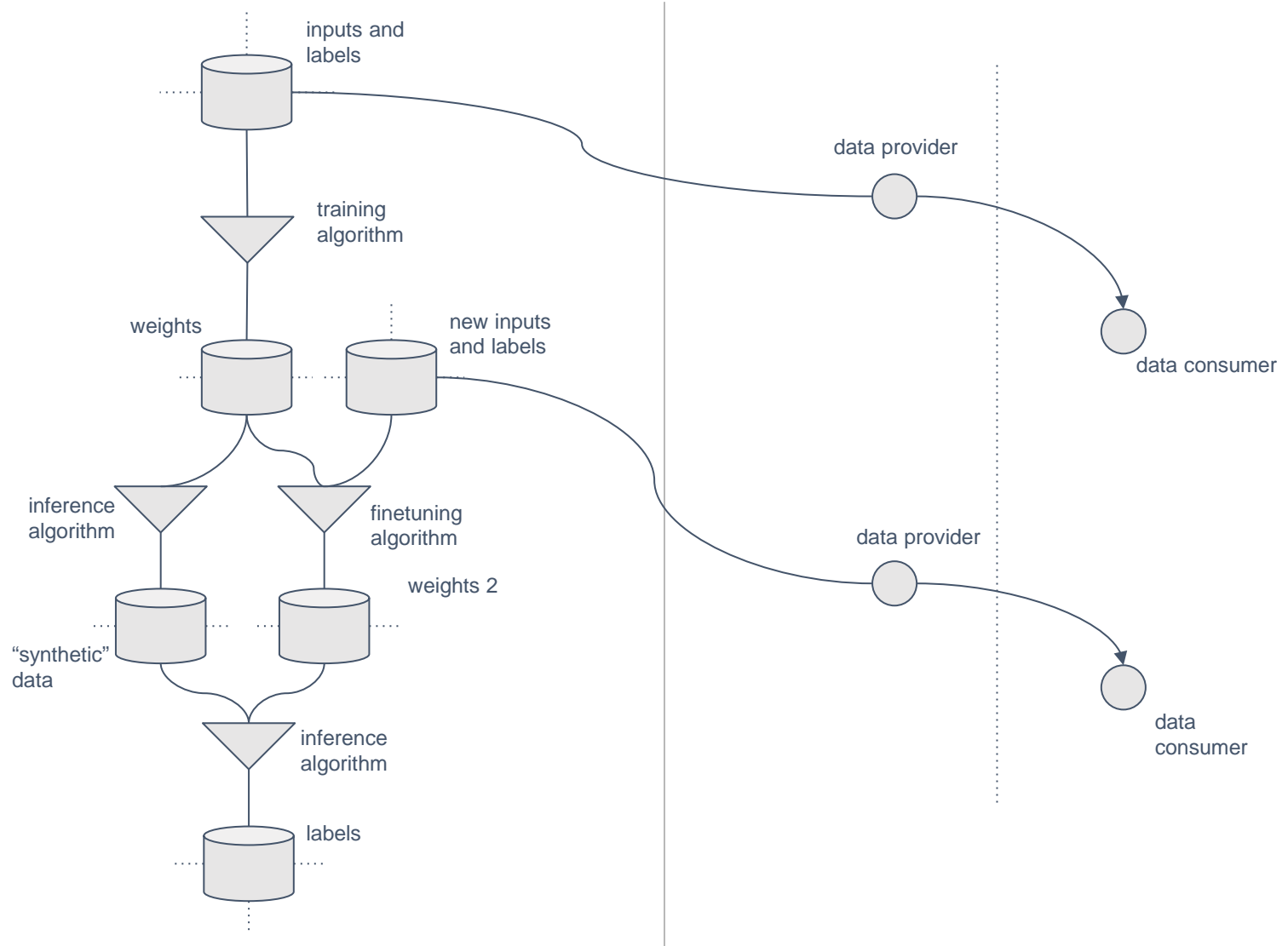
Challenges



Data value chain

Transaction graph

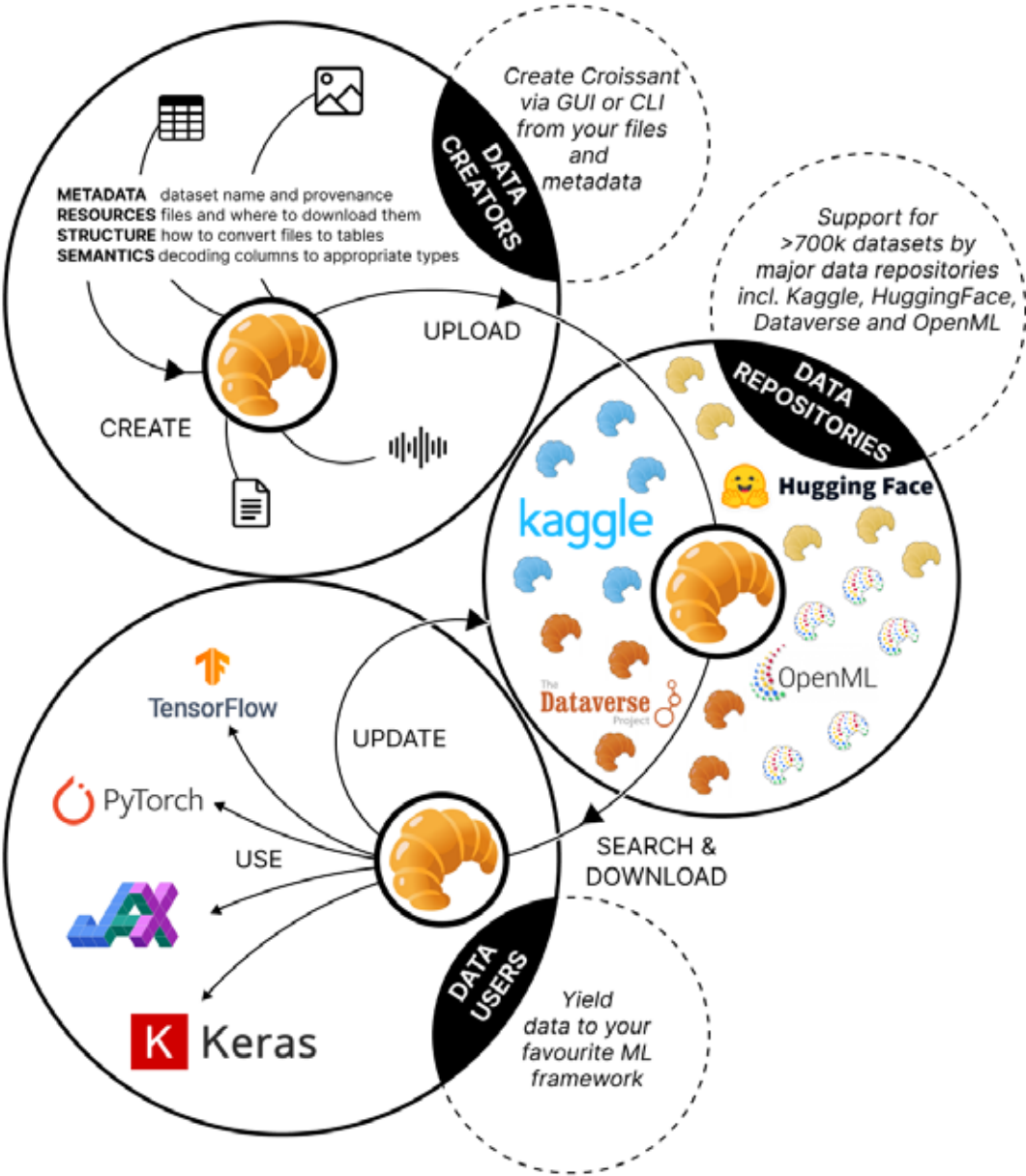
Actor spectrum



- data consumer
 - frontier labs
 - agent wrappers
 - gen pop agents
 - gen pop eng.
- data provider
 - “mom and pop”
 - waking giants
 - web 2 brokers
 - scraping pros

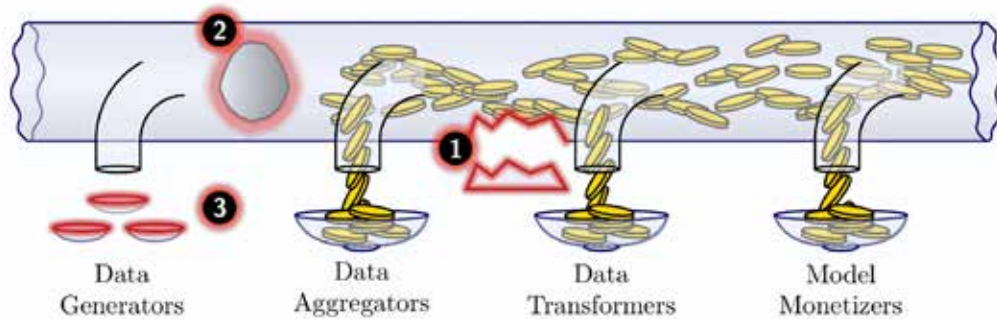
Recent steps

Croissant



Data value chains

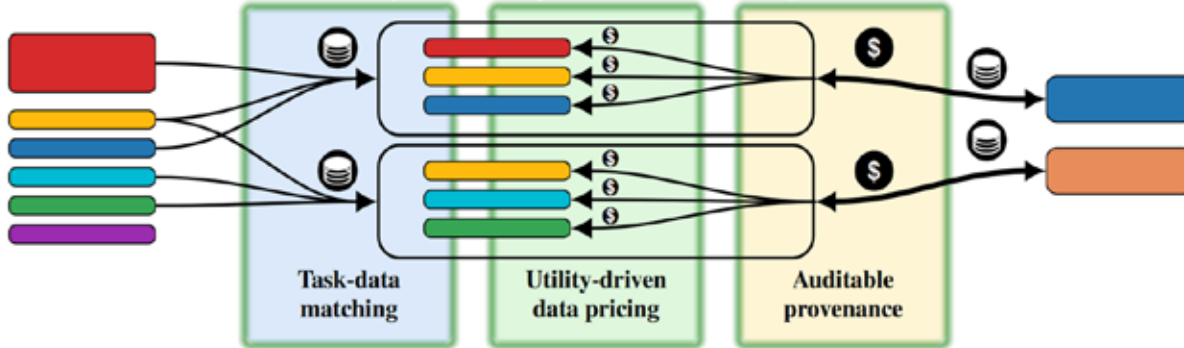
- 1 Invisible Provenance
- 2 Asymmetric Bargaining



- 3 Inefficient Price Discovery

Equitable Data-Value Exchange (EDVEX) Framework

Dynamic task-optimized data bargaining



AI Data Deals

Interactive explorer

A Sustainable AI Economy Needs Data Deals That Work for Generators

Authors: Ruoxi Jia, Luis Ojeda, Wenjie Xiong, Sushu Ge, Jachen T. Wang, Feiyang Kang, Dawn Song

Logos: UNIVERSITY OF CALIFORNIA, Berkeley, UNIVERSITY OF CALIFORNIA, Berkeley

Deal Network

Disclaimer: Our findings compare publicly disclosed deals and public things. However, many transactions are private or under NDA, so our dataset likely undercounts and our classification necessarily simplifies heterogeneous contracts. While we provide a transparent table, completeness cannot be assumed. If you know additional public information of data deals in AI please add them below.

Select nodes to filter the deal table.

Year: 2016 - 2023

Value (million USD): Unlimited - 300

Content Type: Academic, UGC, Images, News, Health, Audio

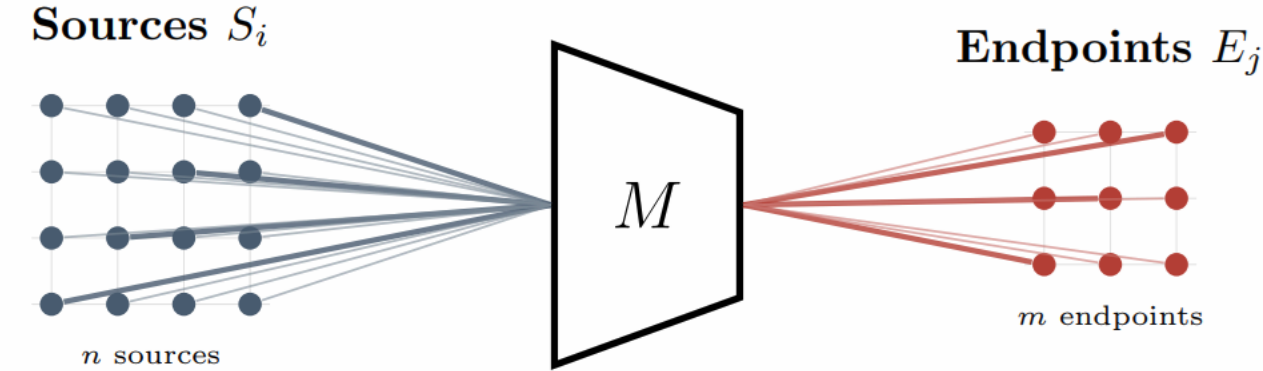
Deal Codes: C, L, S, R, U

Deals (73)

DATA RECEIVER	DATA AGGREGATOR	DATE	TYPE	VALUE	CODES	SOURCE	
Un disclosed	DataDeals AI (Zedjed)	2025	Images	Un disclosed	C	www.zedjed.ai	\$4M
Un disclosed	De Gruyter Brill	2025	Academic	Un disclosed	U	www.degruyter.com	\$4M
Profitable	AAAS	2025	Academic	Un disclosed	C	www.aaas.org	\$4M
Profitable	Research Users	2025	Images	Un disclosed	C	www.research.com	\$4M



Frontier: Closing the loop via multiplexers



— higher utility
— lower utility

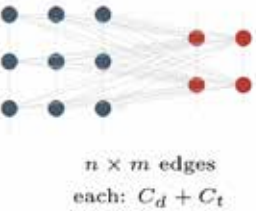
● source $\rightarrow M$
● $M \rightarrow$ endpoint

C_d data cost
 C_t transaction cost

$$\min(C_d + C_t) \text{ s.t. } \sum U \geq \theta$$

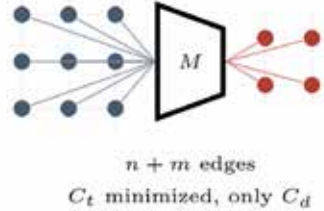
$$n \times m \rightarrow n + m$$

Without Adapter



VS

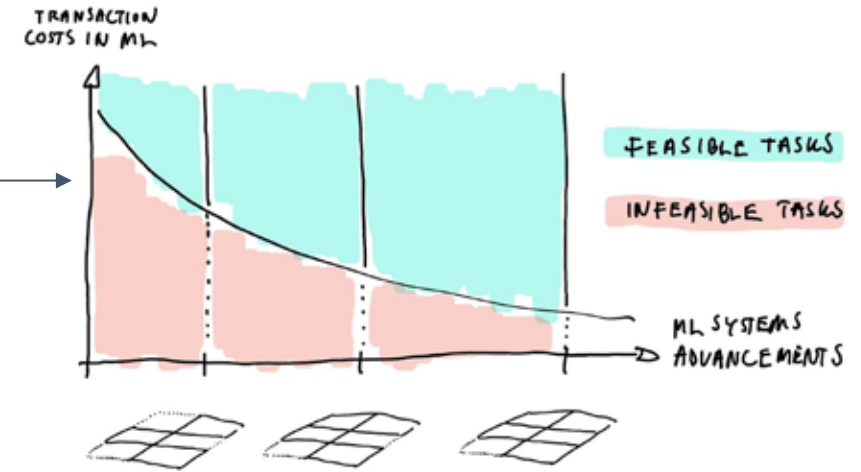
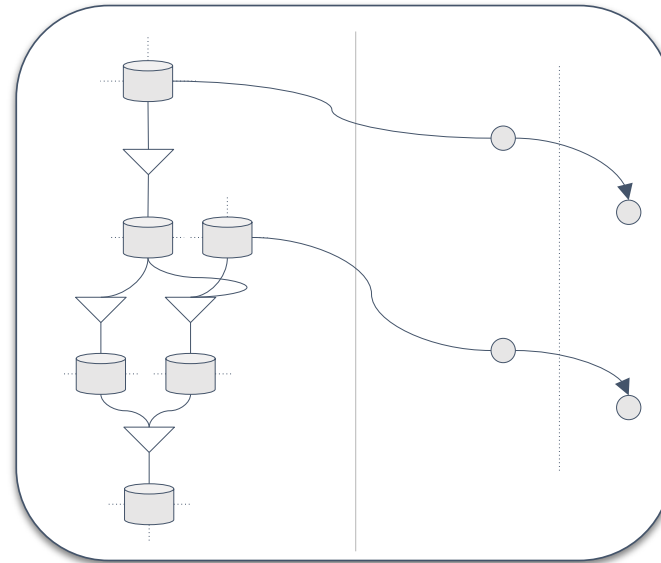
With Adapter



The data multiplexer M acts as a universal adapter between sources S_i and endpoints E_j , routing flows that maximize utility U while minimizing total cost $C = C_d + C_t$. Edge thickness indicates utility; color indicates direction (blue: source $\rightarrow M$, red: $M \rightarrow$ endpoint).



Our north star: A road network for data without congestion



luis@brickroadapp.com
<https://github.com/mlcommons/croissant>

What's next



Expanding the federated catalogue

Onboarding partner institutions and data sources across GI-AI4H benchmarking challenges



More targeted dataset discovery

Using UMLS vocabulary and AI tools to enable detailed search across the catalogue



Tooling for Croissant adoption

Automated metadata generation, validation, and visual editing for data providers



Integration with evaluation pipeline

Connecting catalogue directly to benchmarking challenges for seamless dataset selection



Scaling the data mesh

Decentralized data ownership with cross-functional domain teams managing data as products

→ Now let me hand over to our expert from "BioCroissant" working group for a deep dive into the standard, its specification, tooling, and how you can start using it.

Croissant

A Metadata Format for ML-Ready
Datasets

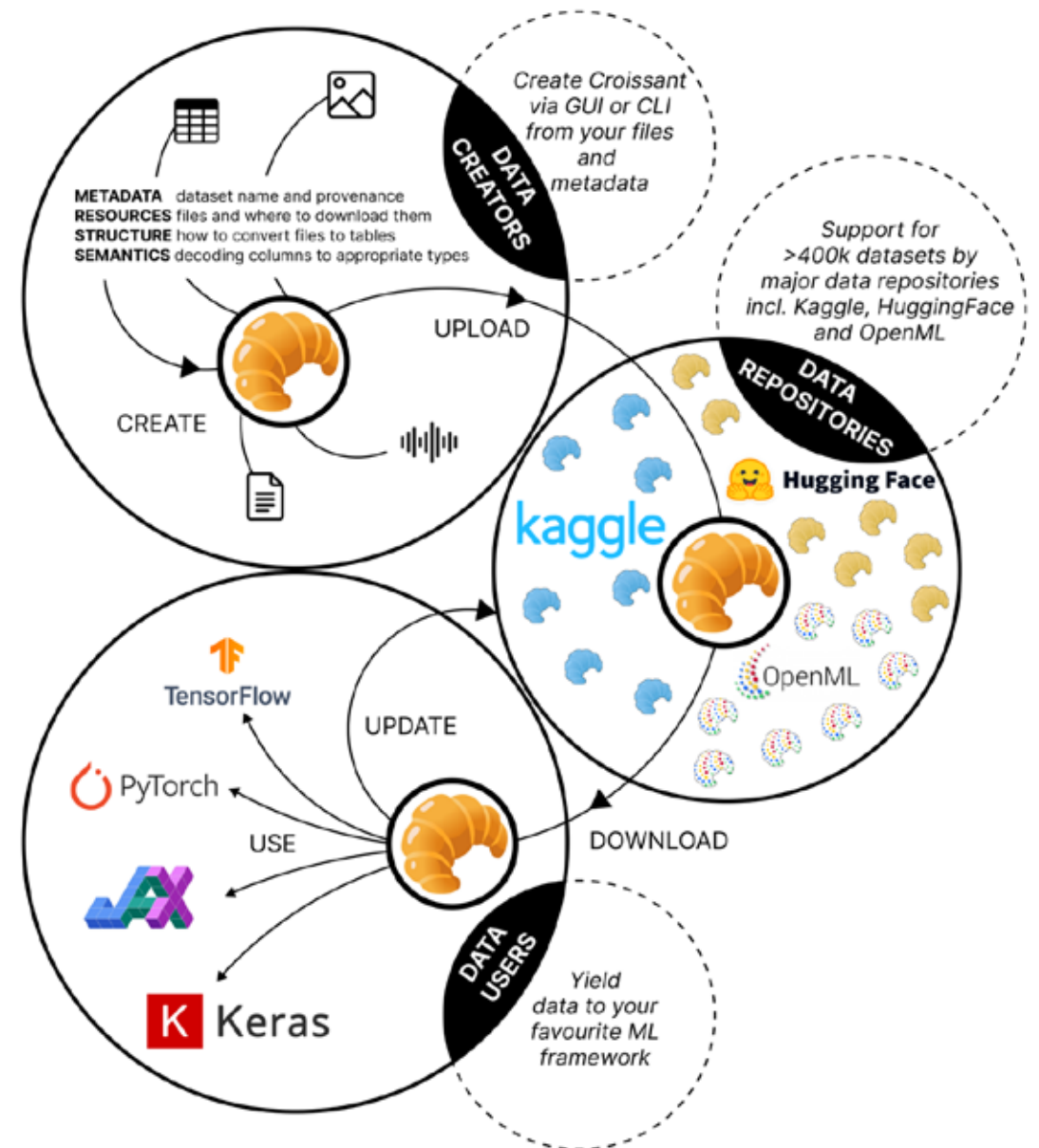
Raeesa Yousaf

Helmholtz Center Munich
German Diabetes Research Centre



Data is Source Code for AI

- **AI quality scales with data quality - but health data lives in silos**
 - Manual assembly does not scale
- **Croissant makes datasets machine-readable without changing the data**
 - One metadata layer across storage locations
 - Compatible for Agentic tools
- **It describes everything AI needs to know about data**
 - Access, structure, licensing, provenance, annotation, and intended use



Croissant on the inside



```
{
  "@type": "sc:Dataset",
  "name": "minimal_example_with_recommended_fields",
  "description": "This is a minimal example, including the required and the recommended fields.",
  "license": "https://creativecommons.org/licenses/by/4.0/",
  "url": "https://example.com/dataset/recipes/minimal-recommended",
  "distribution": [
```

```
{
  "@type": "cr:FileObject",
  "@id": "minimal.csv",
  "name": "minimal.csv",
  "contentType": "data/minimal.csv",
  "encodingFormat": "text/csv",
  "sha256":
"48a7c257f3c90b2a3e529ddd2cca8f4f1bd8e49ed244ef53927649504ac55354"
}
],
```

```
"recordSet": [
{
  "@type": "cr:RecordSet",
  "name": "examples",
  "description": "Records extracted from the example table, with their schema.",
  "field": [
    {
      "@type": "cr:Field",
      "name": "name",
      "description": "The first column contains the name.",
      "dataType": "sc:Text",
      "references": {
        "fileObject": { "@id": "minimal.csv" },
        "extract": {
          "column": "name"
        }
      }
    }
  ]
}
],
```

```
{
  "@type": "cr:Field",
  "name": "age",
  "description": "The second column contains the age.",
  "dataType": "sc:Integer",
  "references": {
    "fileObject": { "@id": "minimal.csv" },
    "extract": {
      "column": "age"
    }
  }
}
]
}
```

Metadata



Structure



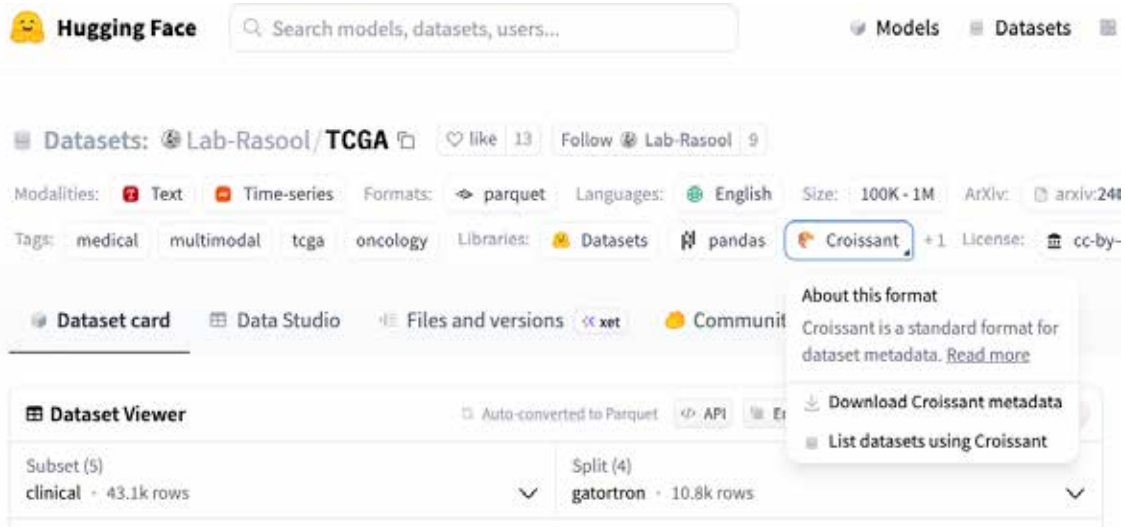
Resources



Semantics



How to use it?



The screenshot shows the Hugging Face Datasets interface for the TCGA dataset. The top navigation bar includes the Hugging Face logo, a search bar, and links for Models and Datasets. The dataset page for 'Lab-Rasool/TCGA' is displayed, with filters for Modalities (Text, Time-series), Formats (parquet), Languages (English), Size (100K-1M), and ArXiv (arxiv:244). Tags include medical, multimodal, tcga, oncology, Libraries (Datasets, pandas, Croissant), and License (cc-by-). The Dataset Viewer shows a subset of 5 clinical rows (43.1k rows) and a split of 4 gatortron rows (10.8k rows). A dropdown menu for the Croissant library is open, showing options like 'About this format', 'Download Croissant metadata', and 'List datasets using Croissant'.

```
pip install mlcroissant
```

```
# 1. Point to a local or remote Croissant file
import mlcroissant as mlc
url = "https://huggingface.co/api/datasets/fashion_mnist/croissant"
# 2. Inspect metadata
print(mlc.Dataset(url).metadata.to_json())
# 3. Use Croissant dataset in your ML workload
import tensorflow_datasets as tfds
builder = tfds.core.dataset_builders.CroissantBuilder(
    jsonld=url,
    record_set_ids=["record_set_fashion_mnist"],
    file_format='array_record',
)
builder.download_and_prepare()
# 4. Split for training/testing
train, test = builder.as_data_source(split=['default[:80%]', 'default[80%:]'])
```

Croissant is a great tool and very easy to use.
But is there something missing to be used in the medical domain?

Why BioCroissant WG is needed – what AI for Health is missing

In analogy to the European Healthcare Data Space (EHDS):

Governance above; technical execution below



Non-normative & voluntary

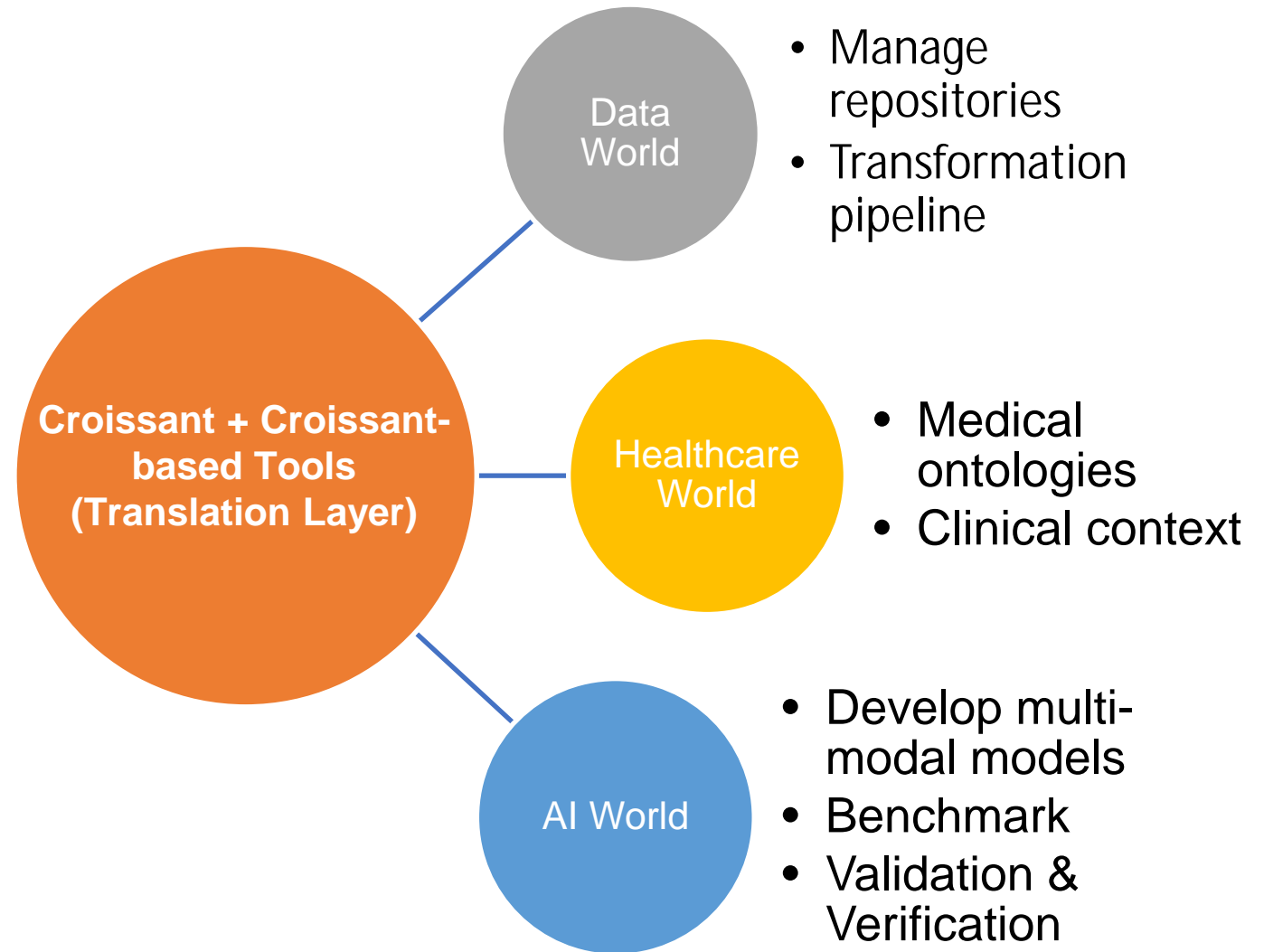
Downstream of EHDS governance

Enables auto-loading of datasets

EHDS-compatible technical profile

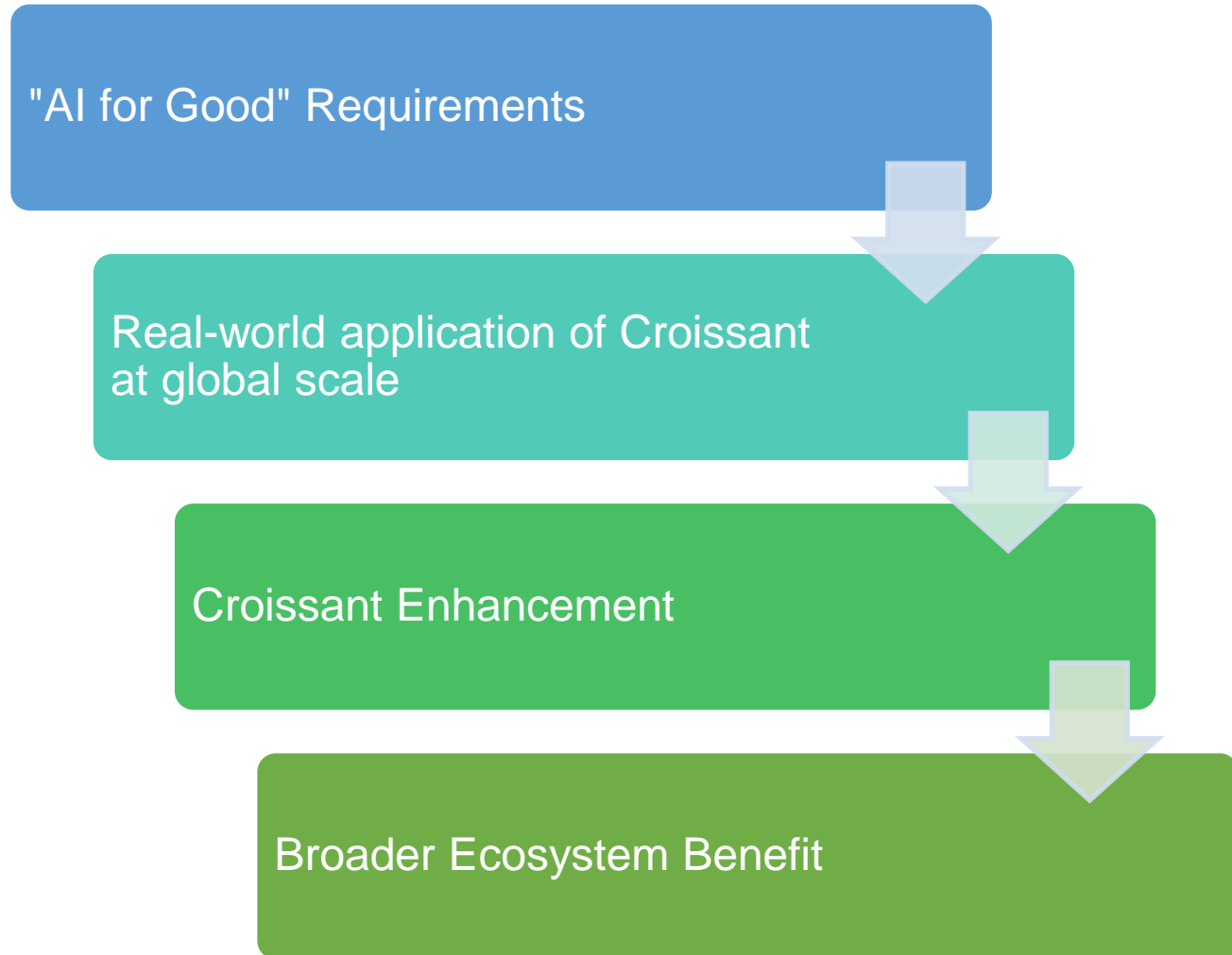
Goal: Bridging Three Communities with one data language

- Smart **metadata wrapping** across modalities, formats, and systems
- **Complexity stays under the hood.**
Tooling handles structure, provenance, licensing, and semantics
- **Experts stay in their lane.**
Each community works in familiar tools and concepts.



From "AI for Health" to Croissant (and Back)

- "AI for Health" brings **real-world requirements** to Croissant
- Global benchmarking turns **theory into practice**
- **Feedback** flows back into Croissant and **improves it**
- Since Croissant is domain-agnostic this "AI for Health" project can indirectly **accelerate standards maturity** for the entire AI field



Call to action: Join the Movement!

Integrate Croissant support into data repositories, analysis tools, and ML frameworks to expand ecosystem

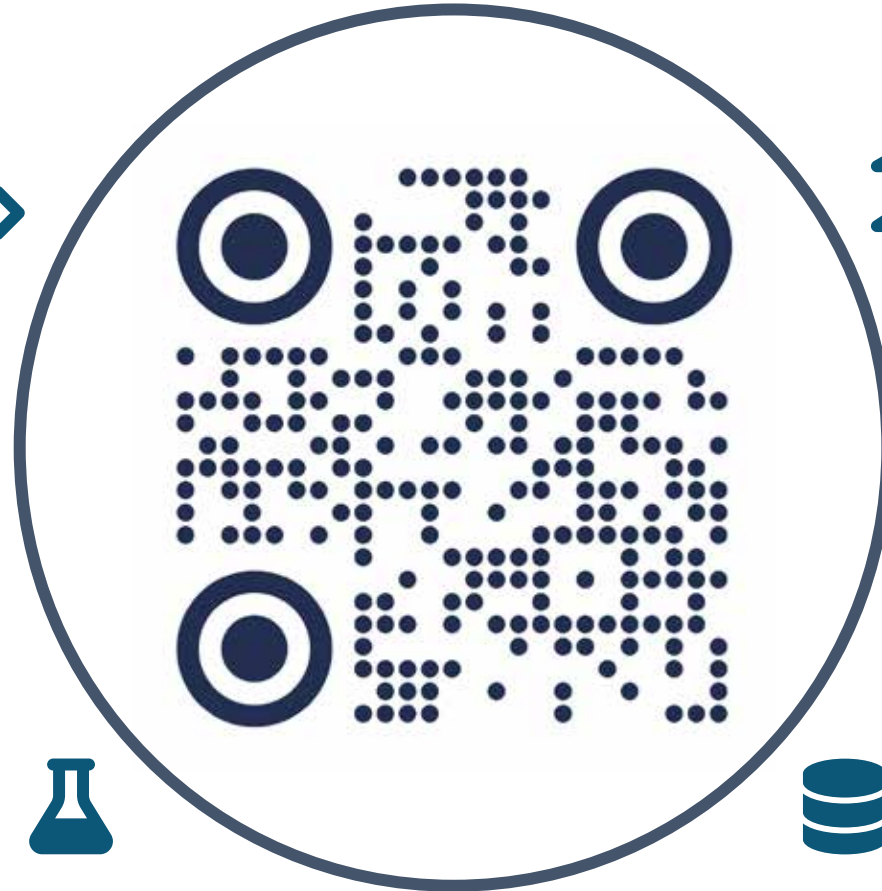
Tool Developers



Reference Croissant standards in data governance frameworks to promote interoperability and FAIR compliance



Policy Makers



Researchers

Contribute to BioCroissant working group, share implementation experiences, and help refine health-specific extensions



Data Providers

Adopt Croissant for health datasets to improve discoverability and enable responsible reuse across platforms

<https://mlcommons.org/working-groups/data/croissant/>

Thank you for your
attention!

