

# Trusted Data Repositories as an Essential Element of the Research Enterprise

---

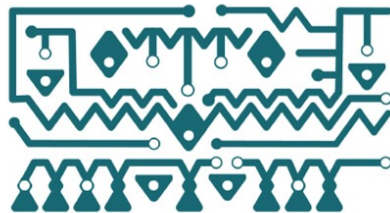
Reyna Broadhurst, Associate Director, ORCID: 0000-0001-6975-6816  
World Data System - International Technology Office, ROR: 01pwmqe95  
2025-Oct-21, eResearch Australasia



# Acknowledgement of Country

We acknowledge the traditional owners of the land on which we meet today - the lands of Turrbal and Yuggera peoples, and pay our respects to their Elders past, present and emerging.

As representatives of the data repository community, we recognize the importance of revising our practices with data to meet the CARE Principles, being mindful that depends on relationships with Indigenous communities and recognizing unique concerns and priorities.



**CARE Principles  
for Indigenous  
Data Governance**

# World Data System (WDS)

The mission of the WDS is to enhance the capabilities, impact, and sustainability of our member data repositories and data services by:



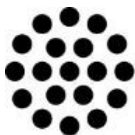
**Creating** trusted communities of scientific data repositories



**Strengthening** the scientific enterprise throughout the entire lifecycle of data and all related components creating first-class data that feeds first-class research output



**Advocating** for accessible data and transparent and reproducible science



**International  
Science Council**

The global voice for science

WDS was initiated by the the predecessor of ISC (ICSU) in 2008, and continues that relationship as an ISC member.



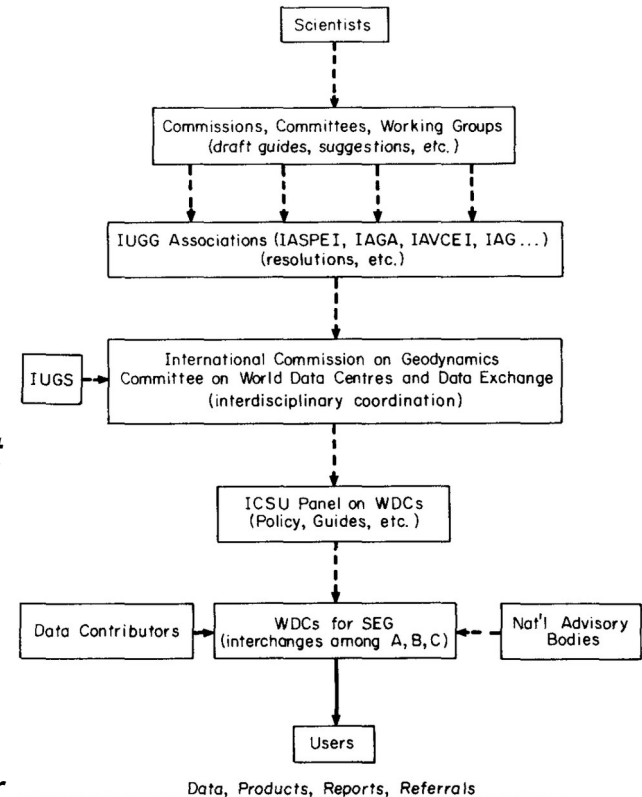
Figure 1. Schematic of Relationships for the World Data Centres for Solid Earth Geophysics

# World Data System Origins

*Data are the lifeblood of science and the sharing of observations and measurements has always been essential to the formulation of new concepts and the evaluation of existing ones.*

*... it was recognized that a more formal mechanism would be needed to provide for the archiving and dissemination of the relatively large quantities of data that would accumulate. To meet this need, the IGY founders established a system of World Data Centers.*

*Modern data centers require a wide range of capabilities including duplicating, microfilming, digitizing, computer graphics, plotting and publishing. They must be able to meet the needs of users for analog or digital data and for summaries, searches, maps and other data products.*



# World Data System History

International Polar Year (IPY)

International Research Council

IPY

International Council for Science

WDS IPO opens in Japan

WDS IPO moves to USA

1899

1931

1957

2007

2018

1882

1919

1932

1998

2008

2021

International Association of Academics

International Council of Scientific Unions

World Data Centers  
Federation  
Astronomical  
Geophysical Analysis  
Services

IPY

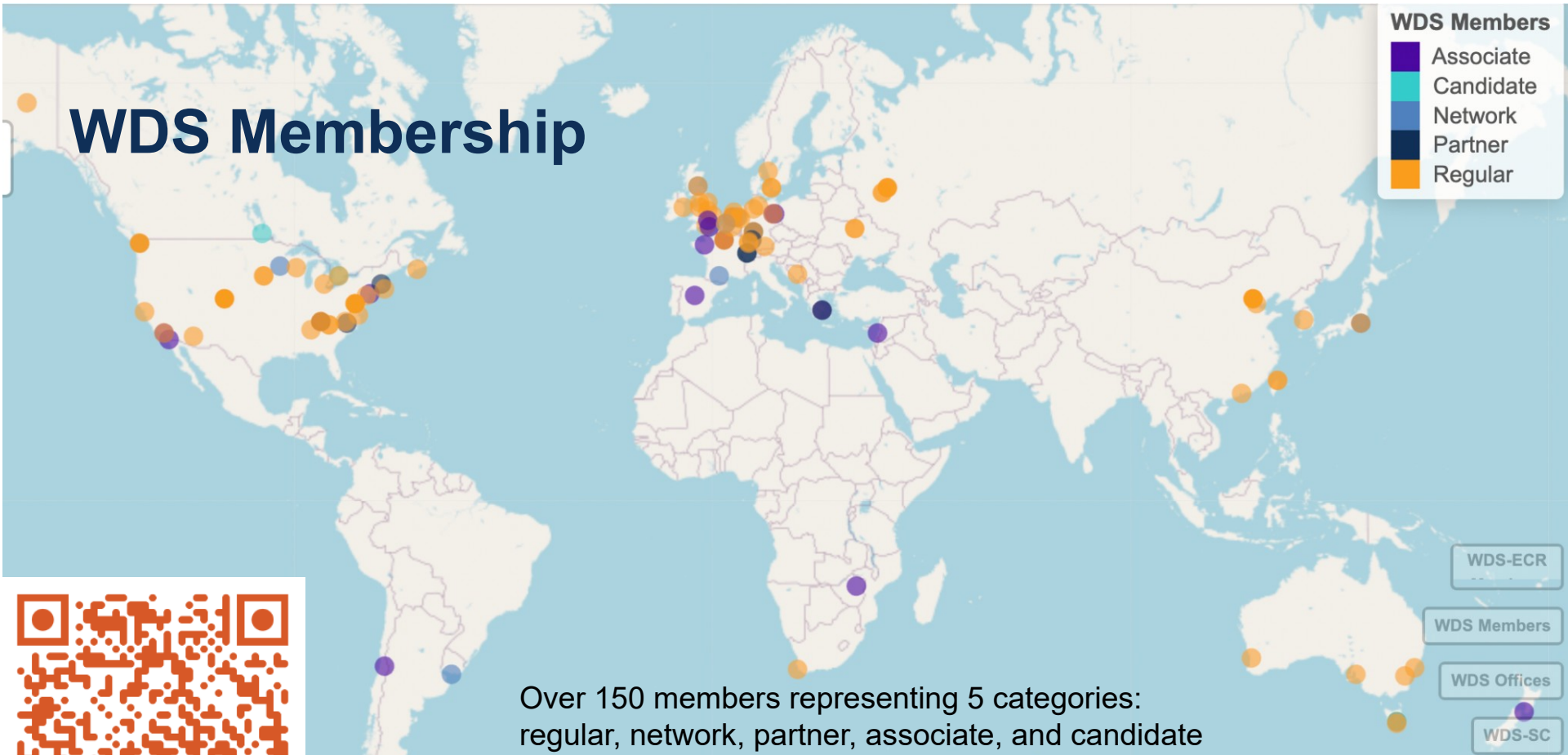
International Science Council

WDS ITO opened

# WDS Membership

## WDS Members

- Associate
- Candidate
- Network
- Partner
- Regular



WDS-ECR

WDS Members

WDS Offices

WDS-SC

Over 150 members representing 5 categories:  
regular, network, partner, associate, and candidate

Become a member:

[https://worlddatasystem.org/members/application\\_membership/](https://worlddatasystem.org/members/application_membership/)



# WDS Scientific Committee

David Castle (Chair)

Maggie Levenstein

Hugh Shanahan (Vice Chair)

Christine Choirat

Juanle Wang

Mamoru Ishii

Yasuhiro Murayama

Devika Madalli

Claudia Medeiros

Dale Peters (Vice Chair)

Indian Ocean

Australia

Libby Liggins

Johnathan Kool

- |                                    |                        |                                      |                        |
|------------------------------------|------------------------|--------------------------------------|------------------------|
| ● Computer Science                 | ● Health Data Science  | ● Knowledge & Information Management | ● Public Policy        |
| ● Data-Intensive Science           | ● High Energy Physics  | ● Marine ecology                     | ● Quantitative Biology |
| ● Geographical Information Systems | ● Informatics & Policy | ● Political & Social Research        | ● Space Weather        |

# WDS Staff

## International Program Office (WDS-IPO)

- PI: Dr. Suzie Allard
- Executive Director: Meredith Goins
- Program Manager: Daniela Santos Oliveira
- Administrative Assistant: Katherine Read
- Informaton Specialist: Samantha Campbell

## International Technology Office (WDS-ITO)

- PI: Dr. David Castle
- Director: Reyna Broadhurst
- Data Science Developers: Andrea Budac, Emilie Altman
- Polar Data Research Associate: Chantelle Verhey



# Garbage In, Garbage Out

Data inputs are crucial to avoid erroneous, biased and harmful outputs.

## Key considerations:

- Gaps & biases due to limitations in coverage (e.g., geographical, demographic, biodiversity, temporal)
- Sensitive/confidential data considerations
- Utilize & report trusted sources of inputs



Prompt: Garbage in, garbage out for AI algorithms  
(leonardo.ai)

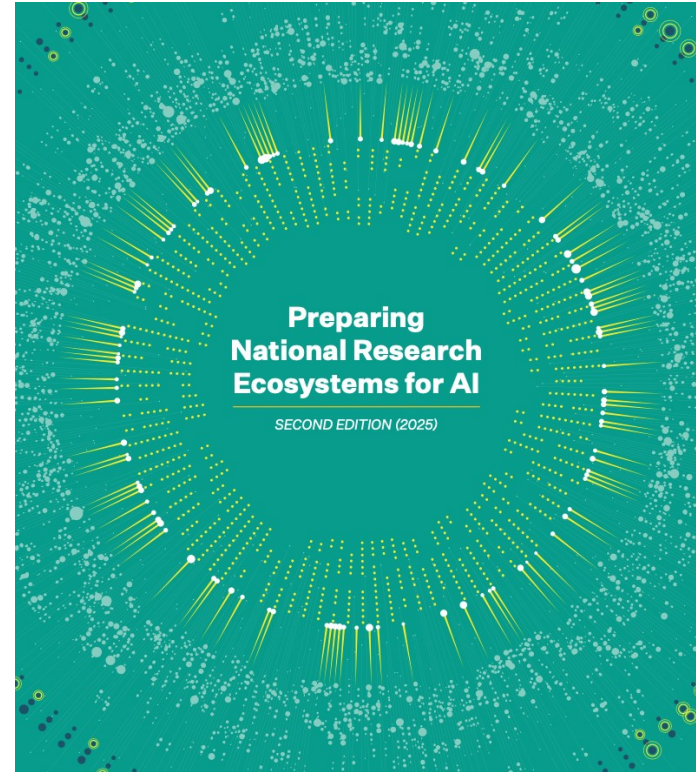
Brooks Hanson, Shelley Stall, Joel Cutcher-Gershenfeld, Kristina Vrouwenvelder, Christopher Wirz, Yuhan (Douglas) Rao & Ge Peng. Garbage in, garbage out: mitigating risks and maximizing benefits of AI in research. *Nature* 623, 28-31 (2023). doi: <https://doi.org/10.1038/d41586-023-03316-8>

## Garbage In, Garbage Out: Value of Curation at Discipline-Based Repositories

“Recognized, quality-assured data sets are particularly needed for generating trust in AI and ML, including through the development of standard training and benchmarking data sets...current data requirements set by funders and journals have inadvertently incentivized researchers to adopt free, quick and easy solutions for preserving their data sets. **Generalist repositories** that instantly register the data set with a digital object identifier (DOI) and generate a supporting web page (landing page) **are increasingly being used....This means that many of the deposited research data and metadata meet only two of the FAIR criteria:** they are findable and accessible. **Interoperability and reusability** require sufficient information about data provenance, calibration, standardization, uncertainties and biases to allow data sets to be combined reliably — which is **especially important for AI-based studies....Disciplinary repositories**, as well as a few generalist ones, **provide this service** — but it takes trained staff and time, usually several weeks at least....We also **urge funders to require that researchers use suitable repositories** as part of their data sharing and management plan. Institutions should support and partner with those, instead of expanding their own generalist repositories.”

# ISC Centre for Science Futures

- AI is a thematic focus area: <https://council.science/home/artificial-intelligence/>
- **Regional workshops** held in Kuala Lumpur, Malaysia (Oct 2023), Santiago de Chile (April 2024) and Oman (Feb 2025)
- Project Lead is David Castle, WDS SC Chair
- Recent report summarizing **literature review and country case studies**: *Preparing National Research Ecosystems for AI*, International Science Council, 2025, <https://council.science/publications/ai-science-systems>,
  - While some commendable progress, there is a general lack of funding, digital infrastructure, regulations, ethics



# Insights from International Science Council Report

Preparing National Research Ecosystems for AI: Strategies and progress in 2024. (n.d.). International Science Council. Retrieved July 23, 2024, from <https://council.science/publications/ai-science-systems/>

4 of the 45 issues that emerged from the literature review were for data quality:

1. **Accuracy**
2. **Bias and exclusion**
3. **Subject orientation of data vs. the interdisciplinary nature of AI research:**  
Most scientific knowledge comes from a specific subject. We need to encode and use it, while enabling communication between domains and allowing for the growing generation of interdisciplinary knowledge.
4. **Data coding and annotation:** AIs, and large language models in particular, require humans to code and annotate the data they use. These individuals must be aware of the risk of embedding cultural differences in the data during the annotation process.

# ESIP Partnership



EARTH SCIENCE  
INFORMATION PARTNERS



**Documentation** enables users to understand the contents of a data set and should provide information and tools to **increase data usability**.

**Quality information** determines the dataset's "**fit-for-purpose**" for AI & affects the trustworthiness of AI applications.

**Access** affects the **efficiency & reproducibility** for AI R&D process.

**Preparation** help identify common **data services and tools** to reduce preprocessing burden for users including AI practitioners.

ESIP Data Readiness Cluster is a forum for community members to:

- Understand users' data needs for AI/ML Research and Development with environmental data
- Develop community standards, leading practices & tools for AI-ready data



# ESIP AI Data Readiness Checklist

## Quality

- Completeness
- Consistency
- **Unbiasedness**
- Timeliness
- **Uncertainty information**
- **Provenance**

## Documentation

- Standard/common metadata
- **Tools & examples (codes, softwares, notebooks)**
- Data dictionary
- Identifiers
- **License information**

## Access

- **Data format** (prefer variety)
- Delivery options (prefer variety, including cloud)
- **Security & privacy**
- **Usage rights**

## Preparation

- **Labels/targets identified**
- **Recommended data splits**
- **Missing values & outliers**
- **Recommended (or not-suitable) for (which techniques would be suitable and which would not)**
- **Past usage (links to documents and/or tutorials)**



# AI Data Readiness Workshop & Pilot Cohort

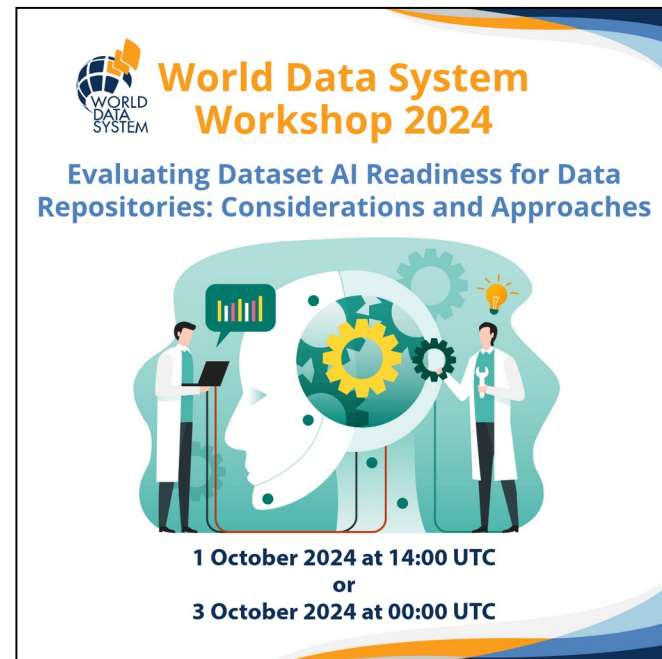
**Target Audience:** data repository representatives who are involved with curating metadata and data (or provisioning of tools for that curation).

## Workshop

- 47 attendees from 12 countries (Australia, Austria, Canada, Chile, Columbia, Finland, Hungary, India, Portugal, South Africa, United Kingdom, USA)
- Disciplines: Agriculture, Earth and Ocean Sciences, Health, Neuroscience, Social Sciences, Generalist, and more

## Pilot

- 12 data repositories included
- Intro, mid-term, and concluding meetings
- Dataset selection, worksheet completion, 1-on-1 sessions
- 2 guest speakers (health data repository experience with checklist, earth sciences expert on bias characterization)



The poster features the World Data System logo at the top left, which consists of a globe with a blue and orange flame-like shape above it. The main title 'World Data System Workshop 2024' is in orange and blue. Below it, the subtitle 'Evaluating Dataset AI Readiness for Data Repositories: Considerations and Approaches' is in blue. The central illustration shows two scientists in white lab coats. One is on the left, looking at a laptop. The other is on the right, holding a test tube. Between them is a large white profile of a human head facing left, with a green gear inside. To the right of the head is another green gear and a yellow lightbulb. A speech bubble with a bar chart is above the first scientist. The background is a light blue and green gradient with abstract shapes. At the bottom, the dates '1 October 2024 at 14:00 UTC or 3 October 2024 at 00:00 UTC' are listed in black.

**World Data System**  
**Workshop 2024**

Evaluating Dataset AI Readiness for Data Repositories: Considerations and Approaches

1 October 2024 at 14:00 UTC  
or  
3 October 2024 at 00:00 UTC

# Cohort Feedback

- Expand **glossary**, especially to define that have multiple interpretations (e.g., raw data, processed data, derived data)
- Provide **rationale** for each question
- Expand drop-down choices (beyond yes, no, not applicable), from **binary to maturity level**
- Include **scoring** and examples
- Add considerations such as **CARE Principles** and whether the dataset has **prior usage by an AI model**
- Consider nuances for different use cases, such as **static data versus dynamic data**

## Further guidance needed:

- Bias
- Quality
- Data versioning
- Dynamic data
- Data curation
- Annotations/labels

**Automation** of checklist is also in consideration, potentially building on open source tools like the F-UJI FAIR dataset checker.

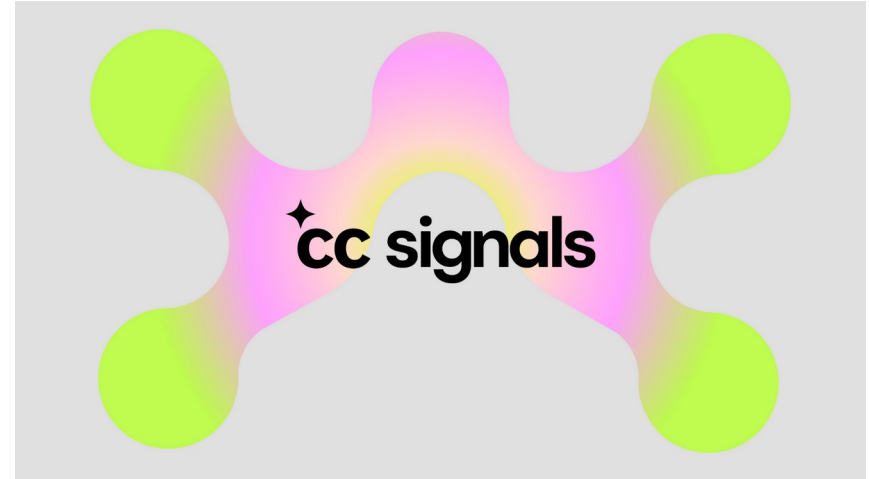
# License Information & Usage Limitations

What is the license for the data?

Is the license standardized and machine-readable?

**CC Signals** - An emerging social contract for AI. The Declaring Party applies CC signals to a set of standard categories: Automated Processing, AI Training, Generative AI Training, AI Use, and Search.

Content-Usage: train-  
genai=n;exceptions=cc-cr



CC Signals © 2025 by Creative Commons is licensed under [CC BY 4.0](#)

# Indigenous Data Governance



Local  
Contexts

Traditional Knowledge (TK) & Biocultural (BC) labels enable customized expression of provenance, protocols, & permissions

## Biological Diversity COP16 - Article 8(j)

- recognizes Indigenous peoples *“efforts in applying, preserving & maintaining their traditional knowledge, innovation & practices in relation to promoting the conservation & sustainable use of biodiversity.”*
- Includes sharing of benefits from the utilization of genetic resources & digital sequence information on genetic resources, as well as traditional knowledge associated with genetic resources.



### POLICY BRIEF

**Recognizing Indigenous Interests:  
Labeling DSI with Provenance  
Metadata**

Jane Anderson, Maui Hudson, Stephany  
RunningHawk Johnson, KatieLee Riddle

# FORCE11 Data Usage Typologies WG The Future of Research Communications and e-Scholarship

<https://force11.org/group/data-usage-typologies/>

**Intent:** to develop a **common typology of data uses**, along with their description and associated characteristics, so that repositories, publishers, institutions, funders, infrastructure providers and meta-researchers can consistently capture, compare and evaluate the ways in which data is used.

## Deliverables

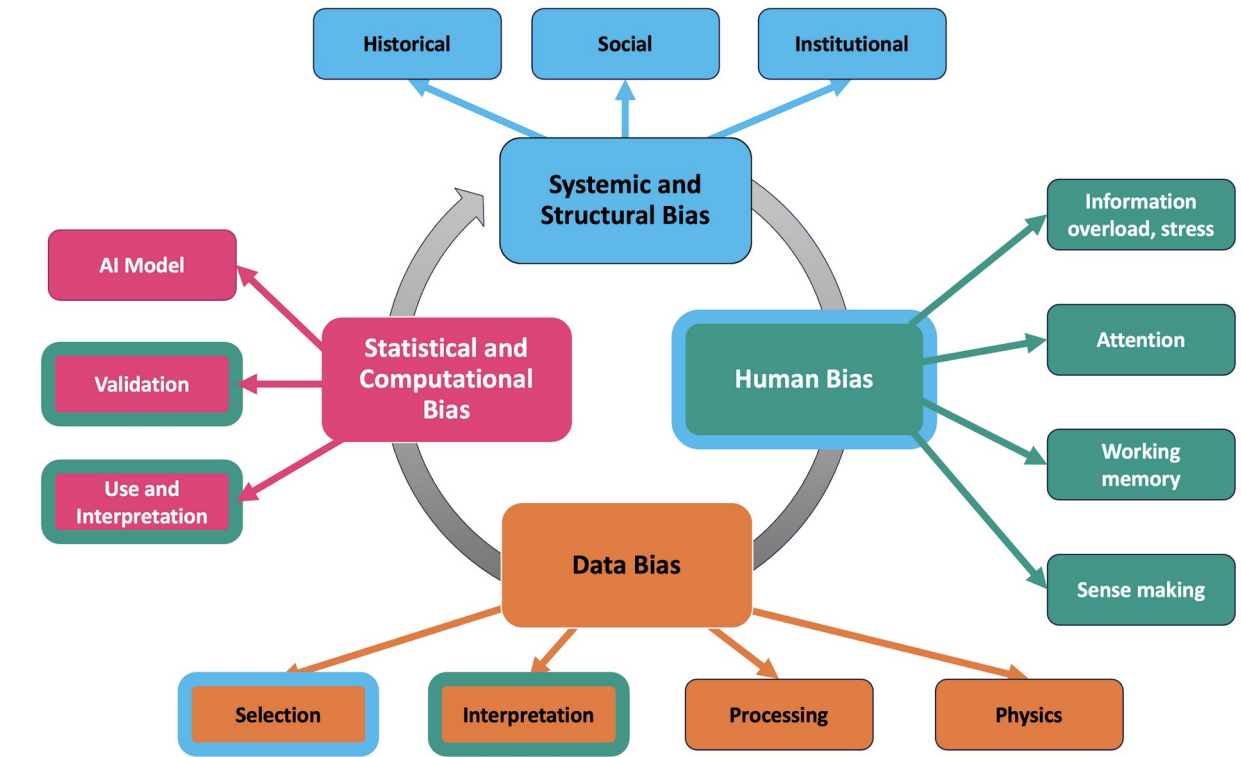
- Codified list of data use types
- Common set of characteristics to identify each of the use types
- Recommendations for how to capture data use types in metadata

While citation in research publications is important, many uses are not captured today:

- Disaster management systems (e.g., early earthquake warning systems, tsunami monitoring systems)
- Weather and climate model data assimilation
- AI models training and benchmarks
- Policy development
- Education
- Technology evaluation (e.g., new instrument testing)
- ...

# Categorization of Bias in AI for Earth Sciences

- Four main categories
  - Systemic & Structural Bias
  - Human Bias
  - Data Bias
  - Statistical and Computational Bias
- First step to measuring and mitigating these biases



McGovern, A., A. Bostrom, M. McGraw, R. J. Chase, D. J. Gagne, I. Ebert-Uphoff, K. D. Musgrave, and A. Schumacher, 2024: Identifying and Categorizing Bias in AI/ML for Earth Sciences. *Bull. Amer. Meteor. Soc.*, **105**, E567–E583, <https://doi.org/10.1175/BAMS-D-23-0196.1>.



# Schema.org Metadata



## [Schema.org Cluster](#)

- Develops domain-specific practices for publishing and harvesting structured Earth science data
- [Dataset metadata guide](#)

## [Ocean Data and Information System \(ODIS\)](#)

- Building an interoperable federated digital ecosystem for all ocean stakeholders, part of the [Ocean Decade vision for a digital representation of the ocean](#)
- Based upon web architectural patterns and use of schema.org/JSON-LD
- [Profile documentation](#)



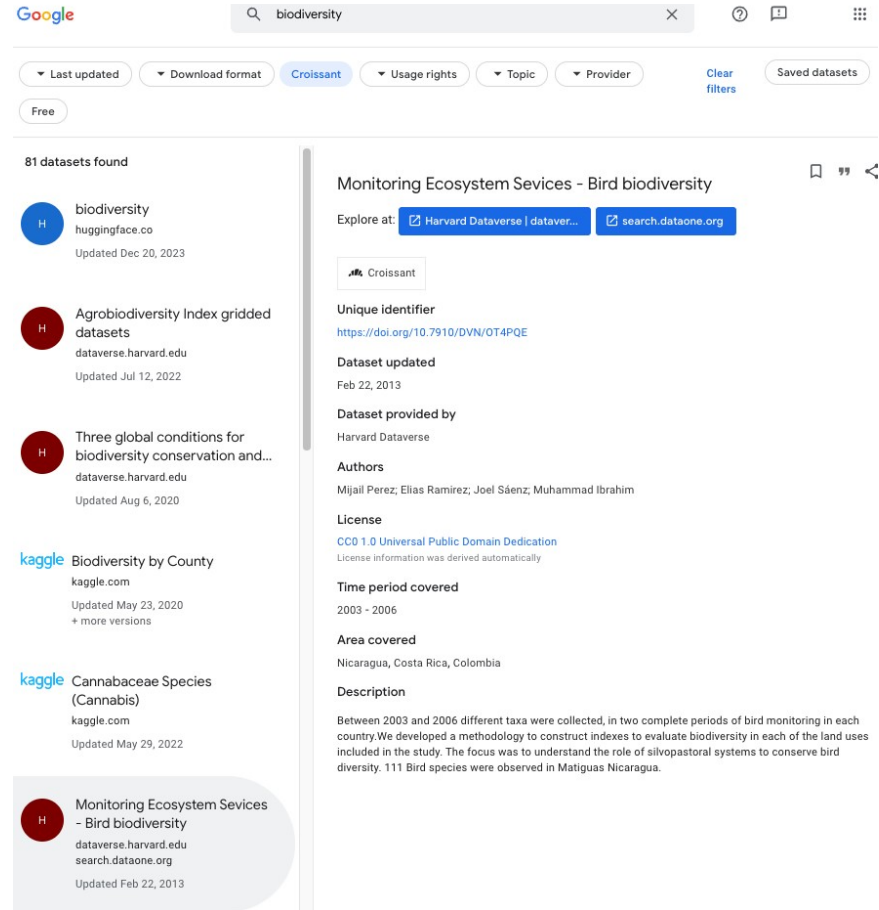
## [Polar Data Search](#)

- Created to aid researchers, stakeholders, & rights holders seeking polar data through a single interface
- Leveraging schema.org metadata for harvesting records from 20 polar data repositories
- [Best practice guide for implementing schema.org](#)



# ML Croissant Extension

- Croissant is a format with 4 layers: metadata, resource file descriptions, data structure, and default ML semantics
- Specification: <https://docs.mlcommons.org/croissant/docs/croissant-spec.html>
- Major outcome of the [ML Commons Croissant Working Group](#)
- Used by Kaggle, Hugging Face, OpenML, etc.
- Croissant filter within <https://datasetsearch.research.google.com/>



The screenshot shows a Google search for 'biodiversity' on the Dataset Search interface. The search results are filtered by 'Croissant' format. The top result is 'Monitoring Ecosystem Services - Bird biodiversity' from Harvard Dataverse, updated on Feb 22, 2013. The interface includes search filters for 'Last updated', 'Download format', 'Croissant', 'Usage rights', 'Topic', and 'Provider'. A 'Free' filter is also visible. The search results list includes:

- biodiversity** (huggingface.co) - Updated Dec 20, 2023
- Agrobiodiversity Index gridded datasets** (dataverse.harvard.edu) - Updated Jul 12, 2022
- Three global conditions for biodiversity conservation and...** (dataverse.harvard.edu) - Updated Aug 6, 2020
- Biodiversity by County** (kaggle.com) - Updated May 23, 2020
- Cannabaceae Species (Cannabis)** (kaggle.com) - Updated May 29, 2022
- Monitoring Ecosystem Services - Bird biodiversity** (dataverse.harvard.edu) - Updated Feb 22, 2013

The detailed view for 'Monitoring Ecosystem Services - Bird biodiversity' shows the following information:

- Unique identifier:** <https://doi.org/10.7910/DVN/OT4PQE>
- Dataset updated:** Feb 22, 2013
- Dataset provided by:** Harvard Dataverse
- Authors:** Mijail Perez, Elias Ramirez, Joel Sáenz, Muhammad Ibrahim
- License:** CC0 1.0 Universal Public Domain Dedication
- Time period covered:** 2003 - 2006
- Area covered:** Nicaragua, Costa Rica, Colombia
- Description:** Between 2003 and 2006 different taxa were collected, in two complete periods of bird monitoring in each country. We developed a methodology to construct indexes to evaluate biodiversity in each of the land uses included in the study. The focus was to understand the role of silvopastoral systems to conserve bird diversity. 111 Bird species were observed in Matiguas Nicaragua.

# Data Spaces, Federated Systems, Data Repositories & Interconnections

Continue to support **data sharing mechanisms and integrations** into data spaces, digital twins and federated systems, regional and disciplinary-specific.]



INTERNATIONAL DATA SPACES ASSOCIATION



Build **relationship metadata and knowledge graphs** between repositories and related entities like harvesters and data spaces with FAIRsharing and re3data



# WDS Data Stewardship Award

This award celebrates **early career individuals** (within 10 years of most recent academic degree) who have significantly improved the quality, integrity and accessibility of research data.

We are now open for applications for the **2026 Data Stewardship Award - due January 16.**

Scan the QR code to learn more criteria, nomination and selection process.

Past award winners: <https://worlddatasystem.org/dsa-previous-winners/>



**World Data System International Program Office (WDS-IPO)** is hosted by the University of Tennessee Oak Ridge Innovation Institute. This work is supported by a cooperative agreement (DE- SC0021915) with the U.S. Department of Energy Office of Science. Contact WDS-IPO at [wds-ipo@utk.edu](mailto:wds-ipo@utk.edu).

**World Data System International Technology Office (WDS-ITO)** is hosted by the University of Victoria. This work is supported by the Digital Research Alliance of Canada and Ocean Networks Canada. Contact WDS-ITO at [ito-webadmin@oceannetworks.ca](mailto:ito-webadmin@oceannetworks.ca).



# Thank You!

Questions?

Visit our website at [worlddatasystem.org](http://worlddatasystem.org) and [wds-ito.org](http://wds-ito.org)

Instagram @wds\_ipo

Facebook @WorldDataSystem.International

LinkedIn @company/world-data-system

Bluesky @wds-ito.org

Vimeo @worlddatasystem

Reyna Broadhurst  
ito-director@oceannetworks.ca

