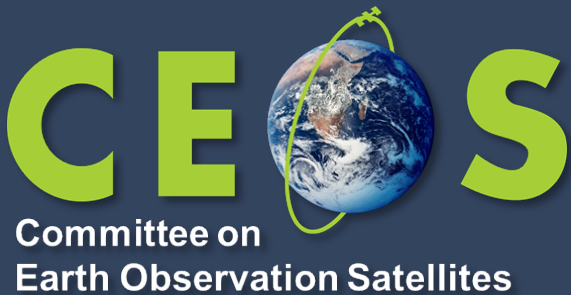# WGISS-54

# FAIR Dataset Quality Information Guidelines

**Dr Ivana Ivánová, Curtin University**

**Agenda ID: 2022.10.03_12.15**

**WGISS-54**

**Tokyo, Japan (JAXA)**

**3-7 October 2022**

CEOS
Committee on
Earth Observation Satellites

# Executive Summary

❖ FAIR DQI community guidelines provide specific advice on ensuring quality metadata compliant with the FAIR principles for the dataset

❖ FAIR DQI community guidelines are a living document developed by international community for international community

❖ Use-Cases on challenges with quality information are wanted!

❖ FAIR DQI guidelines support Priority 3: Support to CEOS Cal/Val Initiatives to increase CEOS Agency Cal/Val Collaboration

# Why FAIR Quality Information?

❖ Increasingly the reuse of a dataset, particularly where multiple datasets are being merged, requires knowledge of the "quality" of the datasets to be merged.

❖ Particularly where datasets are repurposed for use cases beyond what the original creator intended: "quality" information becomes critical.

❖ With the rise of Artificial Intelligence (AI) and Machine Learning (ML), a new interpretation of FAIR is that it stands for "Fully AI Ready": knowing the "quality" of data to be used is essential to avoid erroneous conclusions

**Cloudy, increasingly FAIR; revisiting the FAIR Data guiding principles for the European Open Science Cloud** [Cite]

https://content.iospress.com/articles/information-services-and-use/isu824

Article type: Research Article

Authors: Mons, Barend[a; b; c; *] | Neylon, Cameron[d] | Velterop, Jan[e] | Dumontier, Michel[f] | da Silva Santos, Luiz Olavo Bonino[b; g] | Wilkins...

Affiliations: [a] Leiden Univers... [b] Dutch Techcentre for Life S... Amsterdam, The Netherlands Australia | [e] Independent Op... Institute for Data Science, Ma... Amsterdam, Amsterdam, The... I.N.I.A., Madrid, Spain

Correspondence: [*] Correspon... 9600, 2300 RC Leiden, The Ne...

Keywords: FAIR Data, Open S...

DOI: 10.3233/ISU-170824

Journal: Information Services...

Published: 7 March 2017

⬇ Get PDF

Abstract

The FAIR Data Principles pr... Interoperable, and Reusable... behaviours that researchers... principles should manifest i... Principles has spread, so ha... spread of interpretation, se... the Principles, to clarify bot...

https://danielskatzblog.wordpress.com/2017/06/22/fair-is-not-fair-enough/

June 22, 2017

## FAIR is not fair enough

The FAIR data principles, defined as "a set of guiding principles to make data Findable, Accessible, Interoperable, and Re-usable," came out of a meeting in Jan 2014 that "brought together 25 high level participants representing leading research infrastructures and policy institutes, publishers, semantic web specialists, innovators, computer scientists and experimental (e)Scientists."

The idea of FAIR seems to be catching on, and potentially being applied to other types of objects, such as software. For example, a recent paper, "Four simple recommendations to encourage best practices in research software" (of which I am one of many co-authors), says:

"While the FAIR principles were originally designed for data, they are sufficiently general that their high level concepts can be applied to any digital object including software. Though not all the recommendations from the FAIR data principles directly apply to software, there is good alignment between the OSS recommendations [the software recommendations in the paper] and the FAIR data principles"

eResearch'21 – BoF on FAIR Data Quality Information

❖ Few common quotes:

  ▪ *"We can't use that dataset because it is of poor quality";*

  ▪ *"Don't trust data from sector, organisation or a person: it does not meet OUR quality requirements"*

  ▪ *"Don't trust repository XXXX: their datasets are full of errors and of low quality"*

❖ But when pressed, very few could provide concrete examples of:

  ▪ Exactly what and where the supposed errors were in the dataset;

  ▪ What they were benchmarking the supposedly "poor" quality dataset against

  ▪ None could provide a "community-agreed" reference/best practice document that specified what their expectations on quality were.

❖ "community-agreed" guidelines on quality, preferably at an international level are urgently needed

# An International Effort Came Together...

- ❖ Co-organized by:
  - ▪ ESIP Information Quality Cluster (IQC),
  - ▪ Barcelona Supercomputer Centre Evaluation and Quality Control Team (EQC),
  - ▪ ARDC-supported AU/NZ Data Quality Interest Group (DQIG)
- ❖ 22 International Interdisciplinary Domain Experts:
  - ▪ Data producers (in situ, satellite, model),
  - ▪ Stewards (data/science/technology),
  - ▪ Services providers (data/information/infrastructure),
  - ▪ Data publishers and users
- ❖ from 7 countries (USA, Spain, AU, NZ, Germany, UK, France),
  - ▪ with 22+ affiliations (government, academic, private sectors):
    - ○ Data, science, and service centres, institutional repositories
  - ▪ with expert knowledge from data acquisition or production, data and information management, data publishing, services, and applications.



TOGETHER

WE GO FAR

# Timelines and current status

Initial Discussion (ESIP IQC/BSC EQC)

Pre-ESIP Workshop Announced to Prospective Collaborators

Virtual Pre-ESIP Workshop (July 13, 2020)

Pre-ESIP Workshop Summary and Case Statement
(DOI: 10.31219/osf.io/75b92)

Working Group and Guidelines Development

Public Call-to-Action Statement
(DOI: 10.5334/dsj-2021-019)

Community Review of the Guidelines Document

**Guidelines Document First Baseline**
(DOI: 10.31219/osf.io/xsu4p)

Guidelines Document Maintenance and Update

09/19   02/20   07/20   08/20   09/20   12/20   04/21   10/21

# Four key outputs so far



August 2020 - https://osf.io/75b92/

May 2021 - http://doi.org/10.5334/dsj-2021-019

October 2021 - https://osf.io/xsu4p

March 2022 - http://doi.org/10.5334/dsj-2022-008

# Guidelines development principles

❖ Adapting the FAIR guiding principles (Wilkinson et al. 2016);

❖ Taking a whole dataset-lifecycle approach;

❖ Being quality-attribute and assessment-type agnostic;

❖ Common terminology is essential for enabling interoperability;

❖ Developing for the community by the community:

- Through an iterative process, with continuous engagement with all stakeholders,

- Leveraging the experiences and expertise of a team of interdisciplinary domain experts and community best practices and standards.

# Framework defined by 4 dimensions



Dataset Quality Aspects

**Services** — How well the data has been serviced, supported, and (re)used. Use – Service – Reuse. Quality attributes: Service accessibility, timeliness, security

**Science** — How well the processing algorithm or model has been defined, developed, and validated for intended use. Define – Develop – Validate. Quality attributes: Accuracy, Precision, Uncertainty, Validity

**Stewardship** — How well the data has been curated, preserved, and accessed. Preserve – Document – Access. Quality attributes: Completeness of metadata, Metadata standards, Data accessibility

**Product** — How well the product has been produced, evaluated, and utilized. Produce – Evaluate – Release. Quality attributes: Completeness of data, Data formats, Error sources

Based on: https://doi.org/10.1045/july2017-ramapriyan

# Basic workflow for curating and reporting DQI

**Monitoring & Improvement**

**Quality Specification**
- **Define** and describe the scope of the assessment and associate quality attribute(s) or dimension(s)

**Evaluation Specification**
- **Identify** and describe the assessment method and framework

**Evaluation Execution**
- **Perform** the assessment and capture the results in a structured, human- and machine-readable, standard-based format

**Quality Dissemination**
- **Make** the assessment results readily available and usable to stakeholders and collect feedback for improvement

Based on: Peng et al. (2022). DOI: 10.5334/dsj-2022-008

# FAIR DQI guidelines at a glance

❖ **Guideline 1: Describing Dataset (e.g. version, producer)**
  - Ensure the dataset is findable and accessible.

❖ **Guideline 2: Utilizing a quality assessment model**
  - Ensure the assessment model is findable and accessible.

❖ **Guideline 3: Capturing the assessment method and results**
  - Ensure the quality information is interoperable and reusable (machine end-users).

❖ **Guideline 4: Describing the assessment method, workflow and results**
  - Ensure the quality information is findable, accessible, citable and reusable (human end-users)

❖ **Guideline 5: Reporting the dataset quality information**
  - Ensure the information is FAIR

# FAIR DQI guidelines are really FAIR



Mapping Dataset Quality Information (DQI) Guidelines to FAIR Guiding Principles

❖ At the moment: we are collecting use cases to:

- Ensure that the guidelines are in line with the user communities and their applications;
- Justify the need for best practices in describing quality information to ensure and proper use data;
- Collect examples from multiple application domains on the use of FAIR quality information;
- Provide the community with implementation examples of the guidelines;
- Develop the guidelines for the community by the community;

The template is developed to collect data quality use cases to ensure that the guidelines developed by the international FAIR dataset quality information community gudelines working group (2021) are in line with the user communities and their applications. Please contact Ivana Ivanova at ivana.ivanova@curtin.edu.au for questions regarding the use cases collection effort or Ge Peng at ge.peng@uah.edu for issues with accessing the template.

| ID | Who (Name/Organization) | In what capacity (e.g., data producer; data custodian, funder, ...) | Use-Case Description | Typical object type (e.g., dataset, collection, observation, algorithm, instrument) | Current Data Quality best practice | What quality info is needed in addition to current practice (e.g. license info, provenance info) | What quality indicators make you decide to not use a dataset? | If there is no quality information, what happens? | Additional notes | Contact who can develop the use-case in detail |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Hazard Consortium (500+ orgs) | GIS officer for disaster aftermath recovery | The damage done by Superstorm Sandy in October 2012 was unprecedented in its size and scope. In the aftermath of Sandy, Edison Electric Institute (EEI) members also recognized the need to enhance and formalize the mutual assistance program for national events. In September 2013, EEI's Board of Directors approved a framework to institutionalize the lessons | Spatial datasets | none | Information about trust and reliability of the resource | Data source and producer unknown | The recovery process will be lengthy and therefore more costly (e.g., operational expenses, properties and lives) | | Dave Jones (dave@stormcenter.com) |

❖ Please contribute here.

# FAIR DQI guidelines: path forward

❖ Continue promotion through regular presence at: ESIP, OGC, RDA, SciDataCon, eResearch Australasia…

❖ FAIR DQI guidelines is a living document expected to evolve over time based on user feedback and emerging community best practice

❖ FAIR DQI guidelines are not only for Earth Science datasets – we are expanding the discipline diversity

# Thank you!

ivana.ivanova@curtin.edu.au