

---

---

# Community Guidelines on Documenting and Reporting Dataset Quality

Ge Peng, Sr. Principal Research Scientist

Earth System Science Center/MSFC IMPACT  
The University of Alabama in Huntsville  
ESIP Information Quality Cluster, co-chair

---

---

# Done By

## International FAIR-DQI Community Guidelines Working Group

**Core Writing Team**, Ge Peng, Carlo Lacagnina, Ivana Ivánová, Robert R. Downs,  
Hampapuram K. Ramapriyan, Anette Ganske, Lesley Wyborn,  
Dave Jones, Lucy Bastin, Chung-lin Shie, David F. Moroni

---

---

# Presentation Outline

---

---

- Dataset Quality (Information)
  - FAIR Principles in a nutshell
  - Impact and costs of not sharing
  - Needs and challenges
  - Community guidelines development
  - High-level view of the guidelines
- 
-

**Service Quality Attributes**

**Data Quality Attributes**

**Metadata Quality Attributes**

**Software Quality Attributes**

**Workflow/Procedure  
Quality Attributes**



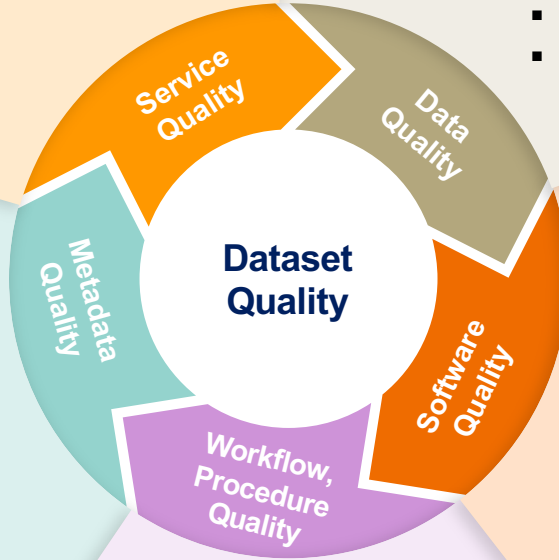
## Service Quality Attributes

## Data Quality Attributes



- Accuracy
- Completeness
- Consistency
- Timeliness
- Integrity
- Validity

## Metadata Quality Attributes



## Software Quality Attributes

## Workflow/Procedure Quality Attributes

## Service Quality Attributes

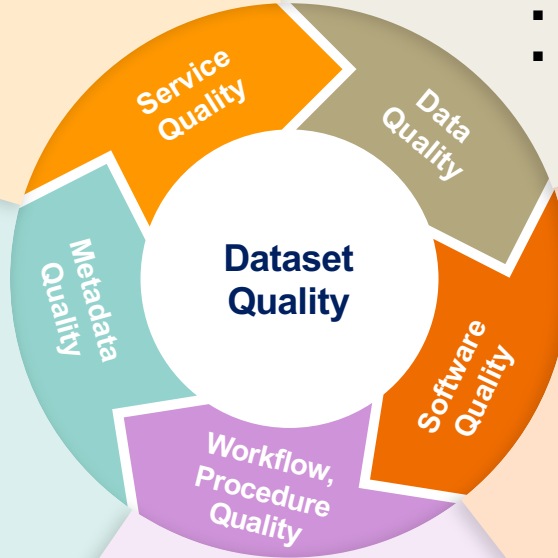
## Data Quality Attributes



- Accuracy
- Completeness
- Consistency
- Timeliness
- Integrity
- Validity

## Metadata Quality Attributes

- Accuracy
- Completeness
- Consistency
- Integrity
- Conformance



## Software Quality Attributes

## Workflow/Procedure Quality Attributes

## Service Quality Attributes

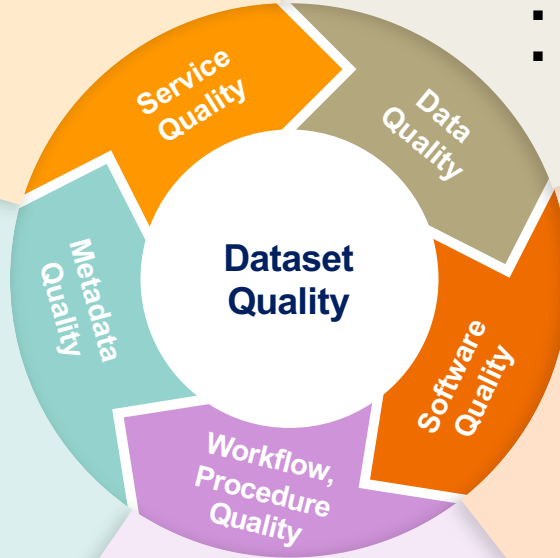
## Data Quality Attributes



- Accuracy
- **Completeness**
- Consistency
- Timeliness
- Integrity
- Validity

## Metadata Quality Attributes

- Accuracy
- **Completeness**
- Consistency
- Integrity
- Conformance



## Software Quality Attributes

## Workflow/Procedure Quality Attributes

## Dataset Quality Information

**Information about quality or the state of**  
data, metadata, software, workflow, procedure, etc.  
through the entire lifecycle of a dataset

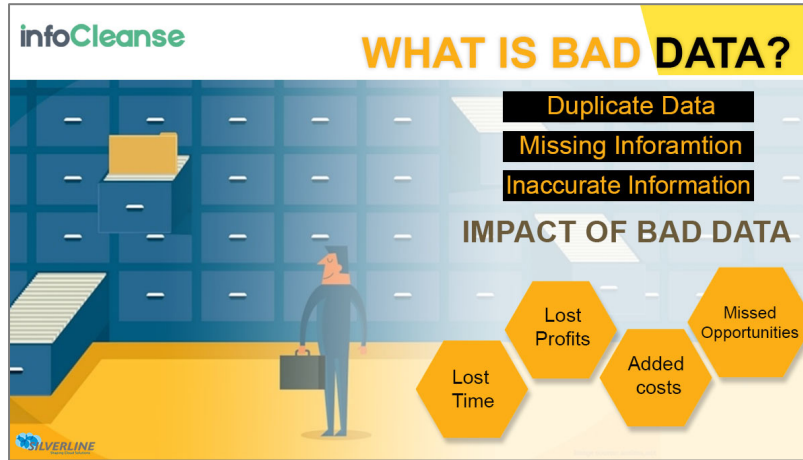
- **Not just information about data quality**
- **Dataset lifecycle approach**

*A dataset* refers to an identifiable collection of data - may contain one or many data files or records in a database in a same data format, having the same variable(s) and product specification(s).



Why Should We Care?

# Impact of 'Bad' Data Quality



Source: <https://www.infocleanse.com/impact-of-poor-data-quality-in-business>



Source: <https://basecapanalytics.com/the-impact-of-bad-data/>

## For Us: Data Producers, Data/Service Providers, Stewards, Data Centers

- Negative in reputation/Trust
- Bad user experiences
- Lost of productivity

**Why you should care about the cleanliness of your data**

**IN DATA WE DON'T TRUST**

**25% of Critical Data in the World's Top Companies is Flawed**

It's no big secret that data entry and errors go hand in hand. Even for Fortune 500 firms, profiling, monitoring and managing the massive amounts of information coming into their businesses is not the easiest of tasks. The worst part, however, is the cost.

**CAUTION**

**Nearly 40% of all company data is found to be inaccurate**

Of 100 companies engaged in a data quality initiative, the best in class found that **23%** of the information they use to make business decisions was inaccurate, while the worst offenders saw an whopping **43%** of their data was just plain bad. Another shocking fact was that

**THE COST OF DIRTY DATA**

140 companies were surveyed and estimated their losses due to inaccurate data:

- For every \$1 million in revenue
- \$5,200,000.00
- On 30 employees
- \$20,000,000.00
- ...the highest 6% of company data over **\$100,000,000.00**

**92%** OF BUSINESSES ADMIT THEIR CONTACT DATA IS NOT ACCURATE!

**66%** OF ORGANIZATIONS BELIEVE THEIR REPUTATION IS AFFECTED BY INACCURATE DATA

**DATA CLEANSING HELPS BUSINESSES**

THE IMPLEMENTATION OF A DATA QUALITY INITIATIVE CAN LEAD TO:

REDUCTIONS			INCREASES	
COSTS, BUDGETS AND OVERHEAD ON MY			ALL THE STUFF YOU WANT MORE OF	
Corporate Budget	IT Budget	Operating Costs	Revenue	Sales
10 - 20%	40 - 50%	40%	15 - 20%	20 - 40%

**THE GROWTH OF INFORMATION**

**Best in class companies claim they can only access 35% of newly added data. Laggards claim 10%**

Companies spend countless amounts of money, time and resources going to incorporate the latest and greatest data sources into their already complex and crowded enterprise management systems, often with less than desirable results. And if you think it's bad now, just wait...

**GROWTH IN GLOBAL DATA**

By 2020, it's estimated that the average organization will have to manage over 30 Zetabytes of data. That means over the next 7 years, data in enterprises will increase by **4400%**

**Halo Business Intelligence** [www.halobi.com](http://www.halobi.com)

Source: <https://www.pinterest.com/ileanagacruz/data-quality>

Why Should We Care?

## Costs of Not Sharing Data

(Collecting Datasets: 10-19% [CrowdFlower 2016])

**For EU: Minimum of €10.2bn per year**

(Source: European Commission and PwC EU Services 2018)

## Costs of Not Sharing Information about Data Quality

(Cleaning and Organizing Data: 60-70% [CrowdFlower 2016])



### Globally

- Productivity lost - redundancy in assessing data quality,
- Million \$ decisions - disaster responses.

# We Need (Consistently Curated) Quality Information

- **Decision-making**

- **Data use**: Informing the reliability and usability of the dataset;
- **Data trust**: Establishing the trust between data providers and consumers, policy-makers.

- **Compliance reporting and open science support**

- Consistently **curated**;
- Readily **available** and **understood** by humans and machines.

- **Support data and information sharing and reuse**

- **Support** new technologies: Interoperable dataset quality information for utilizing Cloud and Machine Learning technologies;
- **Reduce** access barrier: global access and harmonization of quality information.

# Quality Is Complicated!

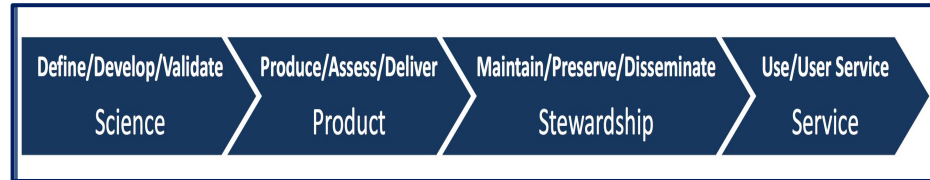
- Multi-dimensionality
- Cross-domain knowledge integration
- Fitness for purpose

# Quality Is Complicated!

- Multi-dimensionality
- Cross-domain knowledge integration
- Fitness for purpose

Quality Attributes	Dimensions
<ul style="list-style-type: none"><li>• accuracy, objectivity, believability, reputation;</li><li>• relevance, timeliness, completeness, value-added, appropriate amount of data;</li><li>• ease of understanding, concise representation and representational consistency, interpretability;</li><li>• accessibility, access security.</li></ul>	<ul style="list-style-type: none"><li>➤ <b>Intrinsic</b></li><li>➤ <b>Contextual</b></li><li>➤ <b>Representational</b></li><li>➤ <b>Accessibility</b></li></ul>

## Dataset lifecycle Stages and Quality Aspects



(Ramapriyan et al. 2017, *D.-Lib Magazine*)

(Wang and Strong 1996, *J. Management Info. Sys.*)

# Quality Is Complicated!

- Multi-dimensionality
- Cross-domain knowledge integration
- Fitness for purpose



# Quality Is Complicated!

- Multi-dimensionality
- Cross-domain knowledge integration
- Fitness for purpose
  - Data user paradigm shift

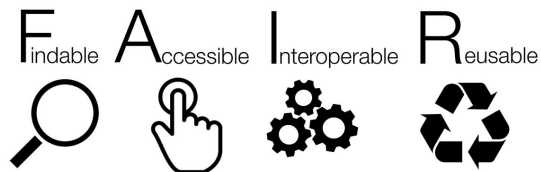


- Making consistently documented quality information readily available

## How To Improve Sharing?

# Adopting FAIR Guiding Principles

(Wilkinson et al. 2016)



(Image by SandyaPundir. CC BY-SA 4.0)



# FAIR Data Guiding Principles

(In a Nutshell)

➤ Uniquely Identifiable  
and Discoverable



## Findable Principle

F1 -> PID

F2 -> Rich Metadata

F3 -> Cross-Reference PIDs

F4 -> Catalog

FINDABLE

ACCESSIBLE

FAIR Data  
Guiding  
Principles

REUSABLE

INTEROPERABLE

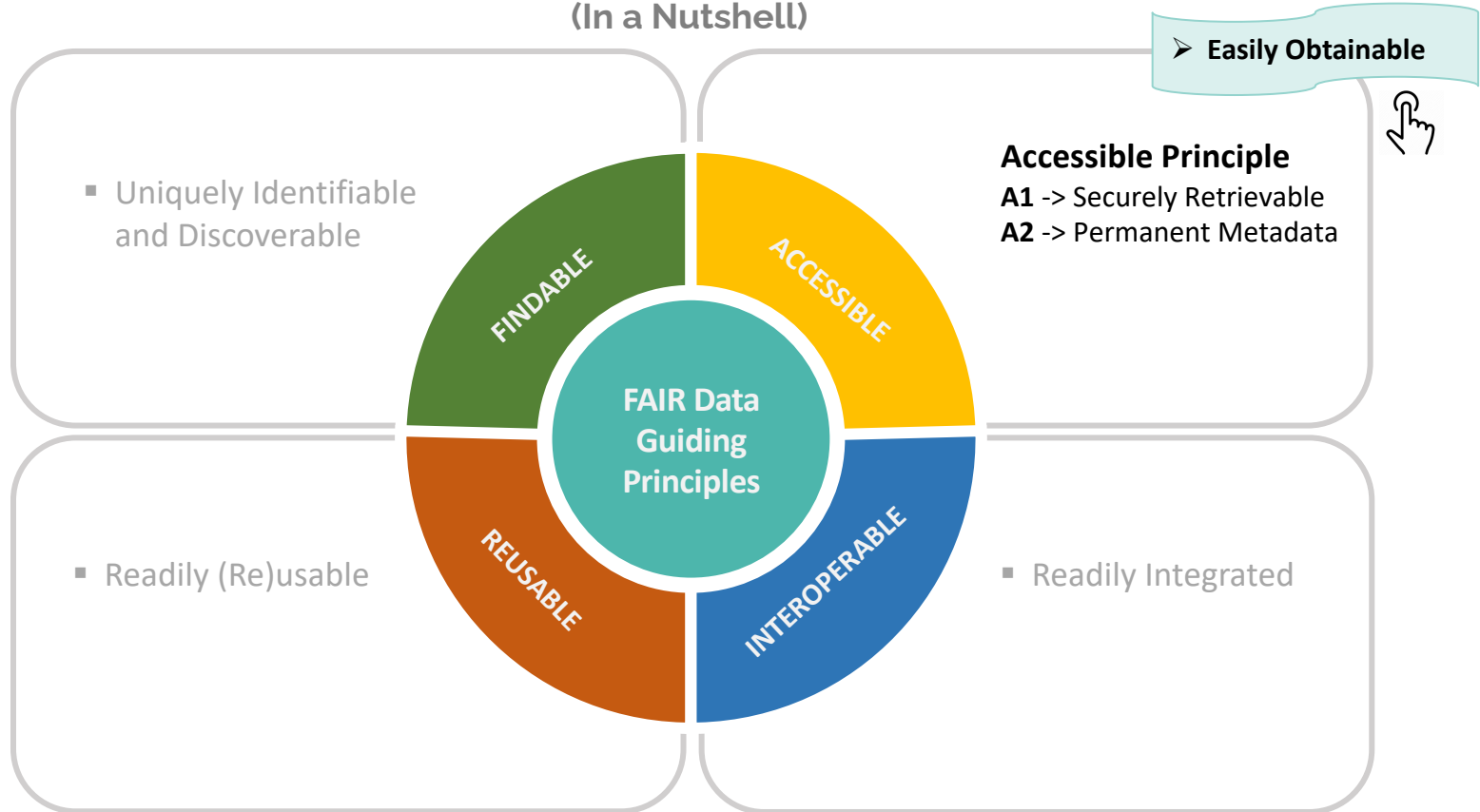
▪ Easily Obtainable

▪ Readily (Re)usable

▪ Readily Integrated

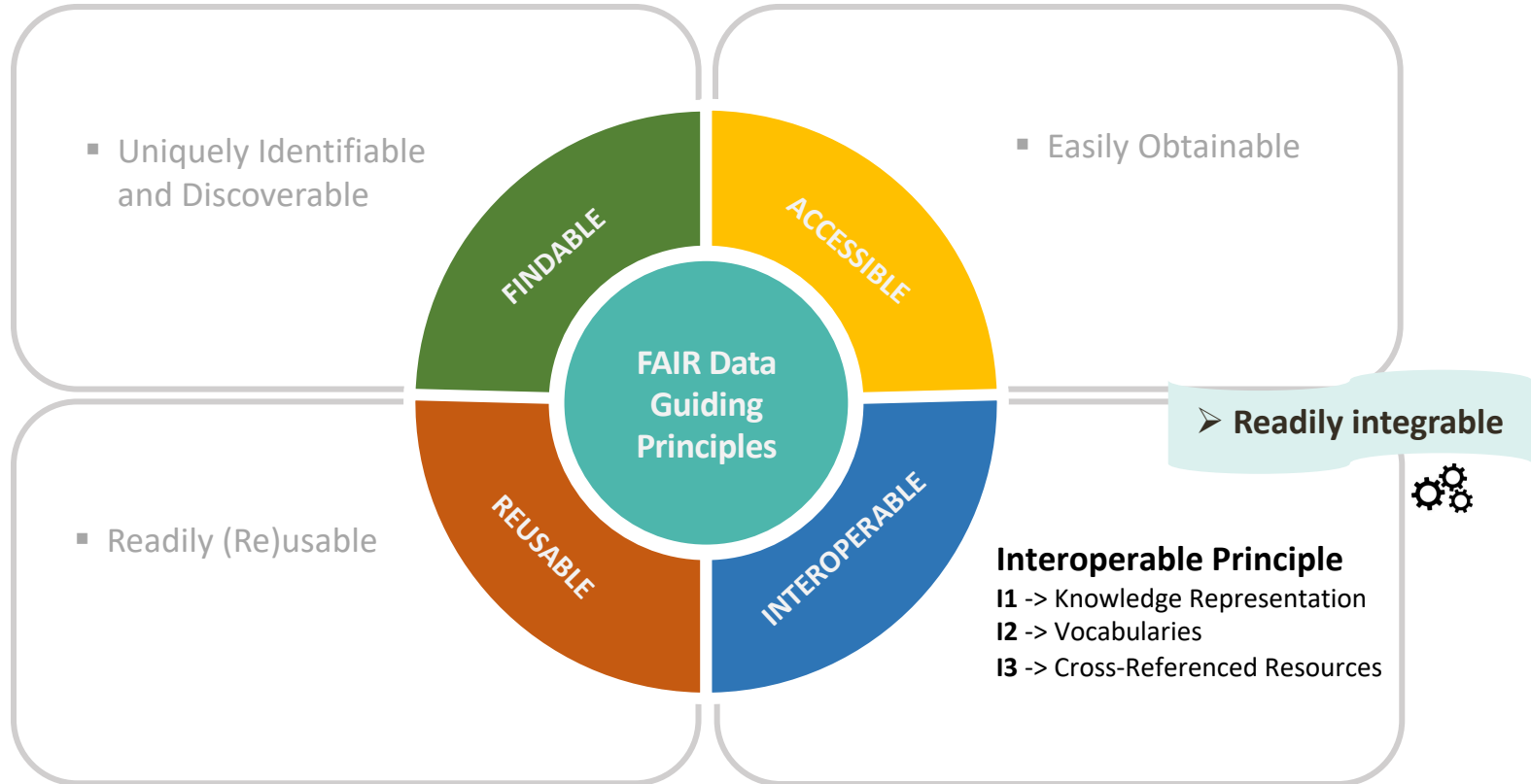
# FAIR Data Guiding Principles

(In a Nutshell)



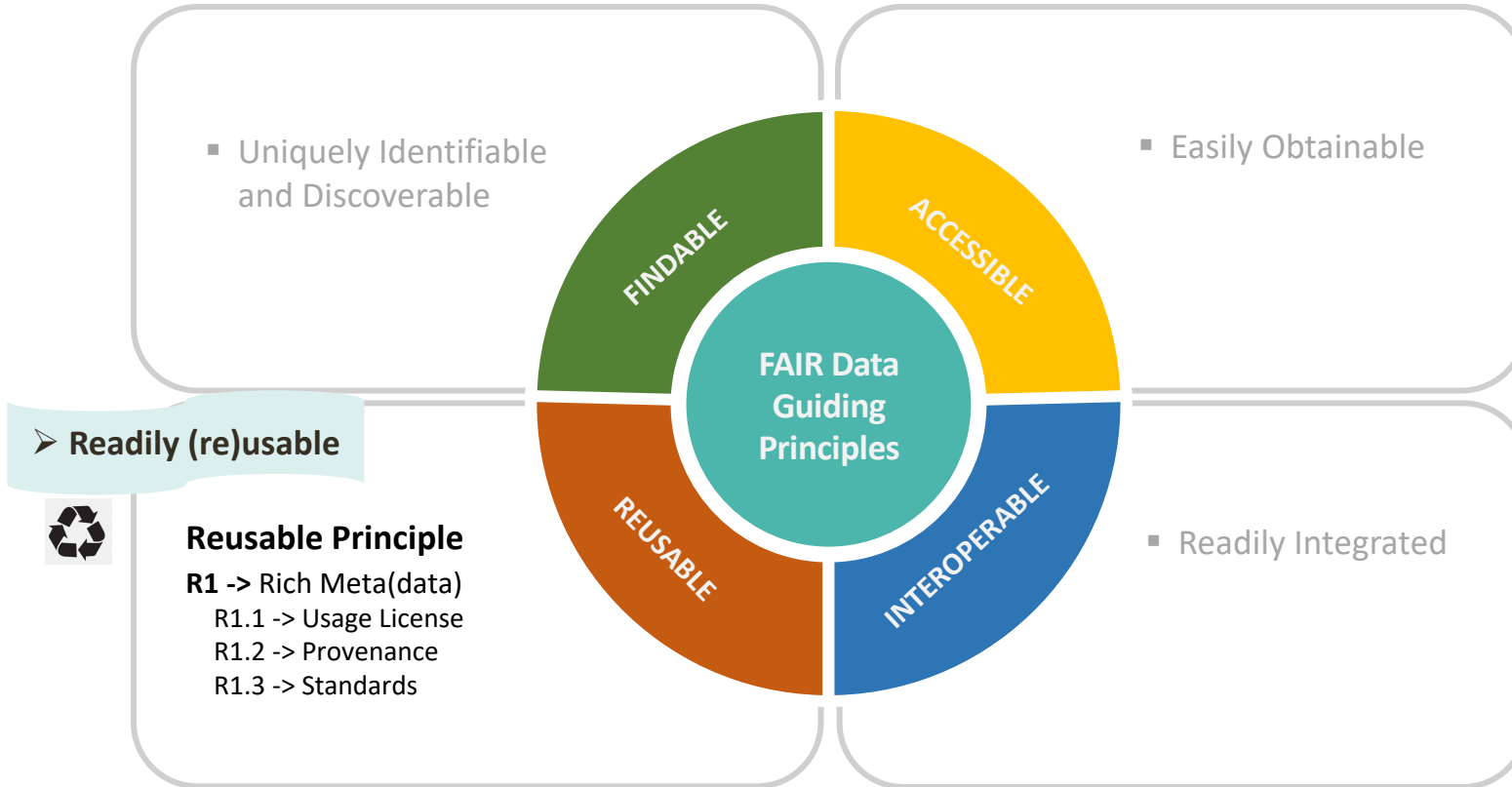
# FAIR Data Guiding Principles

(In a Nutshell)



# FAIR Data Guiding Principles

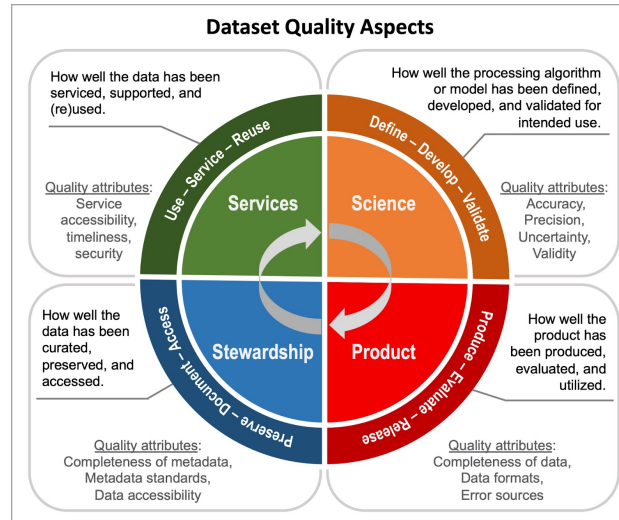
(In a Nutshell)



## How To Consistently Document?

# International FAIR-DQI Community Guidelines

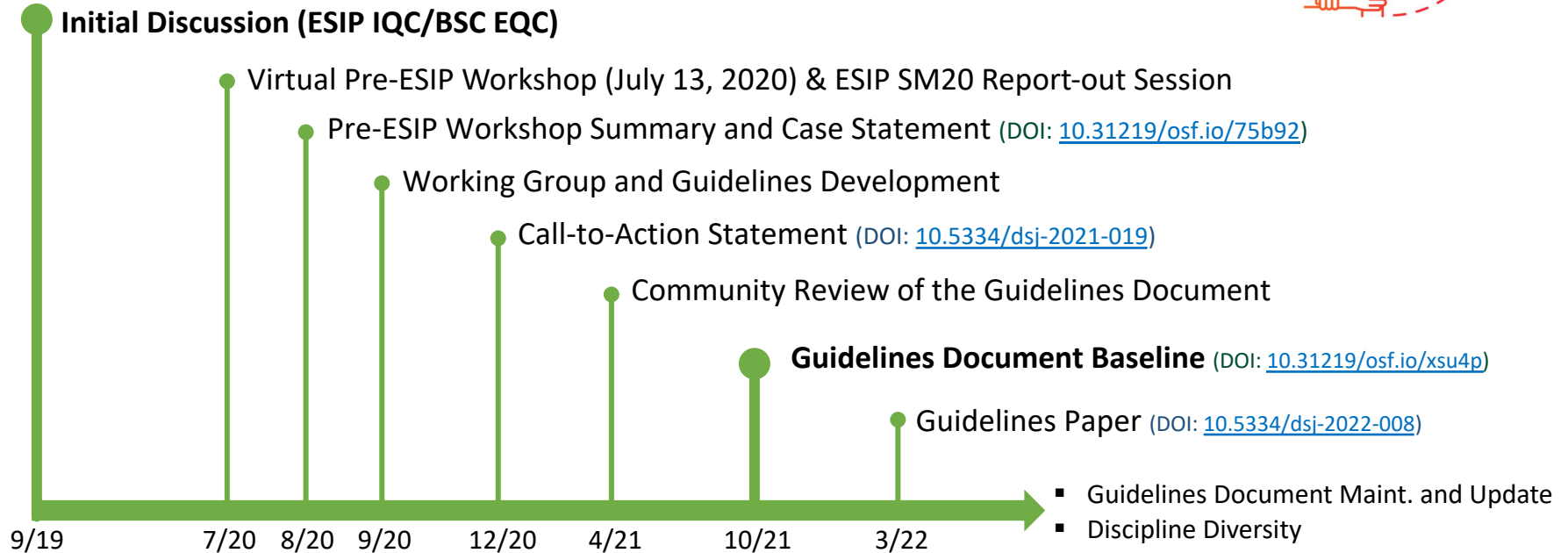
(Peng et al. 2022. DOI: [10.5334/dsj-2022-008](https://doi.org/10.5334/dsj-2022-008))



# Guidelines Development

## Co-organized by

- ESIP Information Quality Cluster (IQC);
- BSC Evaluation and Quality Control (EQC) Team;
- AU/NZ Data Quality Interest Group (DQIG).



# Who We Are

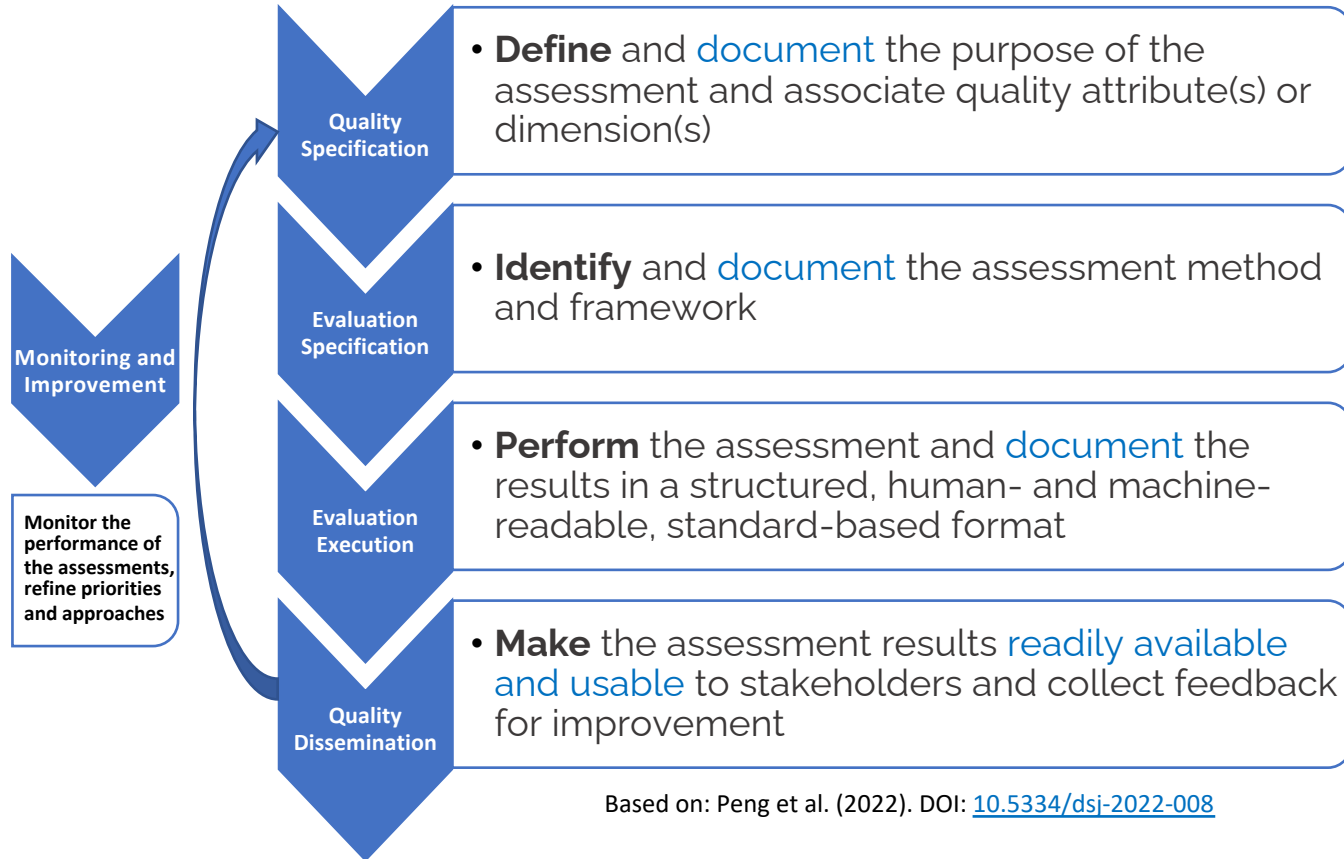
## Community Of Practice

### Participated by

- **22** International Interdisciplinary Domain Experts
  - Data producers (in situ, satellite, model),
  - Stewards (data/science/technology),
  - Services providers (data/information/infrastructure),
  - Data publishers and users.
- **7** Countries (US, ES, AU, NZ, DE, UK, GB)
- **22+** Affiliations (government, academia, private sectors)
  - Data, science and service centers, institutional repositories, companies.
- **Expert Knowledge**
  - Data acquisition or production,
  - Data and information management,
  - Data publishing, services, and
  - Data applications.



# Basic Workflow of Curating and Disseminating DQI



Based on: Peng et al. (2022). DOI: [10.5334/dsj-2022-008](https://doi.org/10.5334/dsj-2022-008)



# FAIR-DQI Guidelines

## (At a Glance)

- **Guideline 1: Describing Dataset**

- Ensure the dataset is findable and accessible

- **Guideline 2: Utilizing a structured quality assessment model**

- Ensure the assessment model is findable and accessible

- **Guideline 3: Documenting the assessment method and results (dataset metadata)**

- Ensure the quality information is interoperable and reusable (*machine end users*)

- **Guideline 4: Documenting the assessment method and results (human-readable document)**

- Ensure the quality information is findable, accessible, citable and reusable (*human end users*)

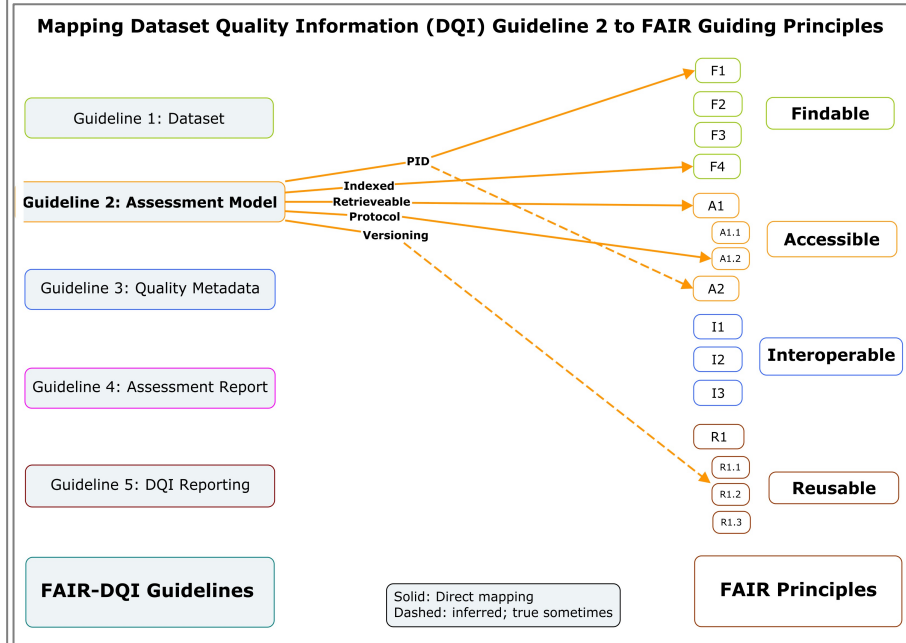
- **Guideline 5: Reporting the dataset quality information**

- Ensure the information is online, findable and readily (re)usable

# Guideline 2 In Detail

**Guideline 2:** Utilize a one or more dimensional, structured quality assessment metric that is:

- 2.1. **versioned** and publicly **available** with a globally unique, persistent and resolvable identifier (**PID**) such as digital object identifier (DOI) and Universally Unique Identifier (UUID);
- 2.2. **registered or indexed** in a searchable resource that supports authentication and authorization, such as Figshare, Zenodo, GitHub, and Dryad; and
- 2.3. **retrievable** by their identifier using an open, free, standardized and universally implementable communications **protocol** such as Hypertext Transfer Protocol Secure (HTTPS) or Open Archives Initiative - Protocol for Metadata Harvesting (OAI-PMH).



- Ensure the assessment model is findable and accessible

Based on: Peng et al. (2022). DOI: [10.5334/dsj-2022-008](https://doi.org/10.5334/dsj-2022-008)

## Examples of dataset quality assessment models and their compliance with Guideline 2

<b>Assessment Model</b>	<b>Scientific Data Stewardship Maturity Matrix</b> (Peng et al. 2015)	<b>Stewardship Maturity Matrix for Climate Data</b> (Peng et al. 2019b)	<b>FAIR Data Maturity Model</b> (RDA FAIR Data Maturity Model Working Group 2020)	<b>Metadata Quality Framework</b> (Bugbee et al. 2021)	<b>Data Quality Analyses and Quality Control Framework</b> (Woo and Gourcuff 2021)
<b>Quality Entity</b> (i.e., attribute, aspect, or dimension)	Stewardship	Stewardship	FAIRness	Metadata	Data
<b>2.1 - Publicly Available</b>	Yes	Yes	Yes	Yes	Yes
<b>2.1 - PID</b>	DOI	DOI	DOI	DOI	DOI
<b>2.2 - Indexed</b>	Data Science Journal	Figshare	Zenodo	Data Science Journal	Integrated Marine Observing System Catalog
<b>2.3 - Retrievable Using free, open, standard-based Protocol</b>	Yes	Yes	Yes	Yes	Yes

# Takeaways

---

---

- Dataset quality is more than just data quality
  - Quality should be considered throughout dataset lifecycle
  - FAIR Principles can help with enhancing sharing of dataset quality information (DQI)
  - FAIR-DQI guidelines can help get started on DQI documentation and sharing
- 
-

**Call-to-action statement** (Peng et al. 2020): [10.5334/dsj-2021-019](https://doi.org/10.5334/dsj-2021-019)

**Guidelines document** (Peng et al. 2021): [10.31219/osf.io/xsu4p](https://doi.org/10.31219/osf.io/xsu4p)

**Guidelines paper** (Peng et al. 2022): [10.5334/dsj-2022-008](https://doi.org/10.5334/dsj-2022-008)

**Join ESIP IQC:** [https://wiki.esipfed.org/Information\\_Quality](https://wiki.esipfed.org/Information_Quality)



# Thank You!

Contact Info: [ge.peng@uah.edu](mailto:ge.peng@uah.edu)

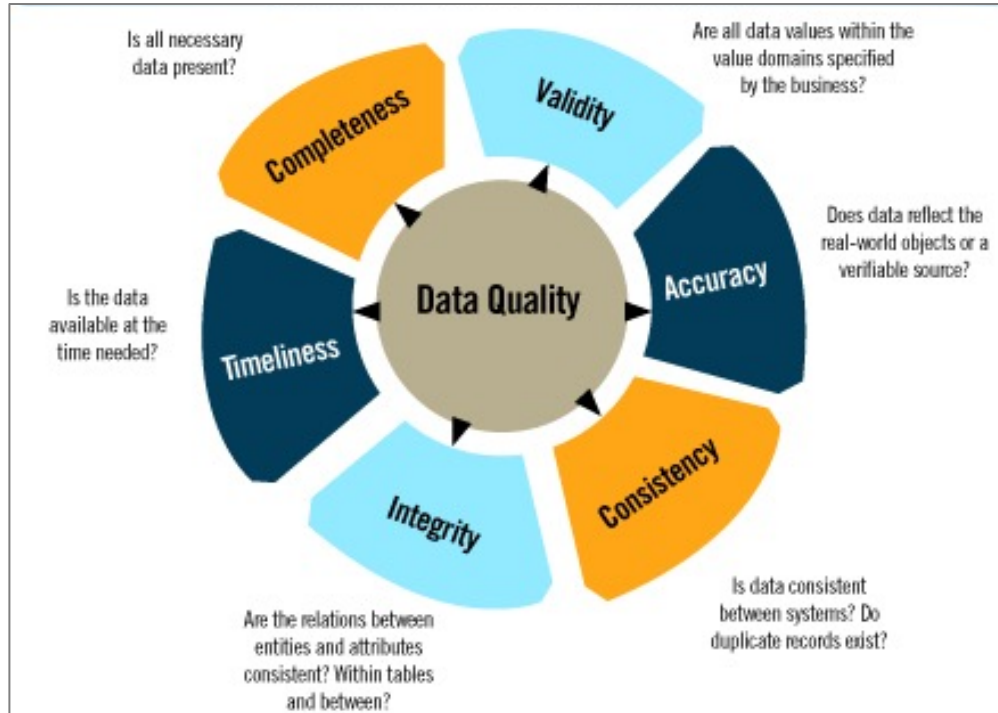
## International FAIR-DQI Community Guidelines Working Group

**Ge Peng, Carlo Lacagnina, Ivana Ivánová,**

Robert R. Downs, Hampapuram Ramapriyan, Anette Ganske,  
Lesley Wyborn, Dave Jones, Lucy Bastin, Chung-Lin Shie,  
David Moroni, Irina Bastrakova, Nancy Ritchey, Mingfang Wu,  
Yaxing Wei, Gilles Larnicol, Kaylin Bugbee, C. Sophie Hou, Ted  
Habermann, Sarah Champion, Gary Berg-Cross, Jeanné le Roux

Ge Peng is supported by NASA Grant **NNM11AA01A** between UAH and MSFC Interagency Implementation and Advanced Concepts Team (IMPACT) project.

# Backup Slides



**Source:** <https://www.kovair.com/blog/data-quality-the-fundamental-of-any-data-migration-project/>



**LEGEND**

- Dataset Lifecycle Stages
- Dataset Quality Aspects
- Document Types
- Metadata Tags
- Metadata Entities

Schematic diagram of dataset lifecycle stages, quality aspects and associated documentation types and metadata tags (MM-\*), and metadata entities. From: Peng et al. (2021). DOI: [10.31219/osf.io/xsu4p](https://doi.org/10.31219/osf.io/xsu4p)

Version: v02r01 202109719  
 POC: gpeng93@gmail.com  
 CC-BY 4.0