

Tools to characterize FAIR and measure FAIRness

Leonardo Guerreiro Azevedo IBM Research – Brazil Iga@br.ibm.com



IBM Research / Oct. 2024 / © 2024 IBM Corporation

Short bio of Leonardo G. Azevedo

- More than 20 years of experience in research and system development
 - Working for industry and government organizations in Databases and Software Engineering
- More than 100 papers published in Brazilian and international conferences and journals
- Research interests: Knowledge Engineering, Data Integration, Semantic Web, FAIR, Service Architecture (MSA and SOA), Business Process Management



BSc. (2000) and PhD. (2005)



Sandwich Doctorate (Jul/2003 to May/2004)



Professor of bachelor and graduation courses (2006 to 2018)



COPPE

MSc. (2001)

Staff Research Scientist (2013 until nowadays)



System Analyst and **Project Manager** (1999 to 2007)

Agenda

- I. Introduction
- -Scenario
- Motivation
- -The FAIR principles
- Main Challenge
- II. Questionnaires for FAIR Characterization
- III. Automated tools for FAIRness assessment
- IV. Discussion

Motivation Scenario



Scientific Research

A researcher wants to compare their dataset resulting from their research with other researchers' datasets

IBM Research / Oct. 2024 / © 2024 IBM Corporation



Find

- Where might the existing dataset have been published?
- How to start the search and using what search tools?
- Which characteristics to filter the datasets should be used?





Access

- Reuse
- Can the dataset be downloaded?
- Does the researcher have permission to use the data? Under what license conditions?
- What are the requirements to integrate the found dataset with the researcher's dataset?
- Are the datasets described with metadata, and metadata in what formats?
- Can the dataset be automatically integrated?

FAIR Principles

Challenge

- Improve data-intensive science for humans and computational agents considering:
 - Discovery
 - Access
 - Integration
 - Analysis

FAIR Data Principles (Wilkinson et al., 2016)

- -15 domain-independent recommendations
- Goal
 - Facilitate data reuse by humans and machines



FAIR at a glance

Keywords

F1. GRUPI F2. rich medata for Findability F3: metadata includes data ID F4: searchable

A1: standard communication protocols A1.1. open, free A1.2. authentication and authorization A2. accessible metadata

I1. knowledge representationI2. vocabulariesI3. qualified references

R1. rich, accurate (meta)data R1.1. clear, accessible license R1.2. detailed provenance R1.3. domain-relevant community standards

IBM Research / Oct. 2024 / © 2024 IBM Corporation

${\bf F} indable$

1. *Identifier* is globally unique, persistent and resolvable (F1)

2. Data has associated metadata for findability (F2)

3. Metadata includes the identifier of its data (F3)

4. Resources are indexed in a searchable manner (F4)

Interoperable

9. Resources use a formal, shared language for *knowledge representation*

- 10. Resources employ vocabularies that are also FAIR
- 11. Resources include qualified *references* to other resources

Accessible

- 5. Identifier uses standardized communications protocols
- 6. This protocol is open, free, and universally implementable

7. This protocol allows for *authentication* and authorization procedures

8. Metadata is reachable, even if its data no longer is

Reusable

- 12. Resources' *attributes* are relevant, accurate and plural
- 13. Resources have a clear data usage *license*
- 14. Resources have detailed *provenance*
- 15. Resources follow domain-relevant community *standards*

FAIR Principles

- FAIR not as in 'just' or 'equitable', but as in 'fitting' or 'suitable'
- "FAIR principles produces digital objects that ensure goals like transparency, reproducibility, and reusability"

Question: How can we assess the adherence of a digital object to the FAIR principles?

However,

- Guidelines are formulated at a high level of abstraction
 - Interpreted and implemented in different ways
 - Contributed to adoption
 - On the other hand
 - Resulted in inconsistent and incompatible interpretations
 - Several works try to harmonize interpretation
 - Papers, e.g., Jacobson et al. (2020)
 - Maturity models, e.g., <u>RDA FAIR Data Maturity Model</u>
 - Metrics, e.g., FAIRsFAIR metrics

Tooling



There are several mechanisms to support the design of FAIR data

- -Guidelines
- -Questionnaires
- -Semi-automated tools
- -Automated tools

Mechanisms' goals

- Characterize digital objects related to the FAIR principles
- -And/or
- Evaluate digital object's FAIRness level

Questionnaires

Essential for

- Overall understanding and
- -Appreciation of the research life cycle

Questionnaire allows for

- -Investigation
- -Acquisition of knowledge about digital objects

Req	The tool should
R1	have open-ended questions
R2	have questions for data and metadata
R3	provide examples that help answer it
R4	deal with all the FAIR principles
R5	allow for a customization
R6	be intended to generate a FAIRness grade
R7	meet indicators or metrics like RDA FAIR Data Maturity Model
R8	support FAIRness test automation
R9	elicit evidence that back the assessment

Existing questionnaires

Source: FAIRassist.org*

- Collects and describes resources to make digital objects FAIR
- -Provided by FAIRsharing
 - a well-known community-driven FAIR initiative

Questionnaires	R1	R2	R3	R4	R5	R6	R7	R8	R9
FAIR Data Self Assessment Tool	0	0			0		Ο	0	0
FAIRDat		0		Ο			0	0	0
FAIRDataBR	Ο	0	Ο		Ο		Ο	Ο	Ο
Data Sterwardship Wizard					Ο	0	0	Ο	Ο
FAIR Implementation Profile (FIP)			0	0	0	0	0	0	0

Labels:

- Requirement is totally supported
- O Requirement is not supported
- equirement is partially supported

The Improved FAIR Characterizaiton Questionnaire



- -Questions
 - For contextualization
 - To characterize digital objects' properties
- Meet questionnaire requirements
 - R1) Open-ended questions
 - R2) Separate questions for data and metadata
 - R3) Present help examples
 - R4) Deal with all the FAIR principles
 - R5) Allow for customization
 - ...

IBM Research / Oct. 2024 / © 2024 IBM Corporation

*Azevedo, L. G., Tesolin, J., Banaggia G., Cerqueira R. "An Improved Questionnaire for FAIR Characterization". In: 3rd Workshop on Metadata and Research (objects) Management for Linked Open Science (DaMaLOS 2023), 2023. DOI: https://repository.publisso.de/resource/frl:6444993

Questions for Findability

-F1: characterize identifer for data and metadata

- What is the main identifier (ID) of the data ...?
- Are there other attributes used to identify the data? If so, what are they?
- Is the data ID globally unique or unique in the dataset domain or for a specific context?
- -F2: characterize the richness of metadata
 - Which metadata schemas, if any, are used to describe the data?
 - What kinds of metadata are used to describe the data?

- -F3: characterize data and metadata linkage
 - What technology does link metadata to the data (and vice-versa)?
 - How are the metadata and data linked?
- -F4: characterize data and meta indexation
 - Which searchable resource is used to register or index the metadata?
 - How is the metadata available or indexed? (E.g., as a static web page, in a database, JSON returned from an API call)

Questions for Accessibility

- Protocol characterization for data and metadata
 - A1: Which communication protocols are used to access the metadata?
 - A1.1: Is the protocol used to access the metadata standardized, open, free, and universally implementable?
 - A1.2: What security mechanisms are used for metadata access, such as those used for authentication and authorization, and access conditions and access levels?

- -A2: characterize data and metadata storage
 - Are data and metadata independently stored?
 - What is the metadata longevity plan?
 - What is the data longevity plan?

Questions for Interoperability

- -I1: characterize knowledge representation and format used for metadata
 - What is the knowledge representation used for metadata? E.g., Relational, Document, Key-Value, Graph ...
 - Is the knowledge representation used for metadata formal, accessible, shared, and broadly applicable?
 - In what format the knowledge representation used for metadata is provided? E.g., eXtensible Markup Language (XML), Turtle (TTL), JSON, JSON-LD ...

- I2: characterize vocabularies for data and metadata
 - Which structured vocabularies are used for metadata?
 - Are these vocabularies used for metadata FAIR in their own right?
- -I3: characterize qualified references
 - Which qualified references does the metadata include to other data or metadata?

Questions for Reusability

- R1: characterizes metadata accuracy and relevant attributes
 - What are the relevant metadata attributes?
 - What is the required accuracy of each metadata attribute, if any?
- -R1.1: characterizes metadata license
 - Which usage license is used for metadata?
 - Is the metadata usage license clear?
 - Is the metadata usage license accessible?

- -R1.2: characterizes provenance
 - Which metadata schemas are used for describing the provenance of the data?
 - What attributes are used for data provenance?
- R1.3: characterizes community standards employed for metadata
 - What are the domain relevant community standards for metadata?
 - Does the metadata under assessment meet these domain-relevant community standards?

Applying the Questionnaire

PubChem

- Open chemistry database at the NIH (National Institutes of Health)
- Chemical information resource for scientists, students, and the general public since 2004
 - Chemical structures
 - Identifiers
 - Chemical and physical properties
 - Biological activities
 - Patents
 - Health
 - Safety
 - Toxicity data
 - ...
- Data sources come from government agencies, chemical vendors, journal publishers etc.



Applying the Questionnaire

- -Source of Information
 - "Perfluorooctanoic acid" compound web page
 - PubChem search web page
 - PubChemRDF
 - Some of PubChem domains represented using semantic web concepts

NIH National Library of Medicine Pub(C)hem Q Search PubChem About Docs Submit Contact COMPOUND SUMMARY 77 Cite Download Perfluorooctanoic acid CONTENTS Title and Summar PubChem CID 9554 1 Structures 2 Names and Identifiers Structure 3 Chemical and Physical Prope 4 Spectral Information 5 Related Records 6 Chamical Vandors 7 Pharmacology and Bioch Chemical Safet 8 Use and Manufacturing 9 Identification 10 Safety and Hazard 11 Toxicity ical Safety Summary (LCSS) Datasheet 12 Associated Disorders and Di Molecular Formula C8HF15O 13 Literature 14 Patents PERFLUOROOCTANOIC ACID Synonyms Pentadecafluorooctanoic acid 15 Interactions and Pathways 335-67-1 16 Biological Test Results PFOA 17 Taxonom

"Perfluorooctanoic acid" at PubChem

PubChemRDF



PubChem Introduction (Documentation)

IBM Research / Oct. 2024 / © 2024 IBM Corporation

Applying the Questionnaire

- -64% of our questionnaire was answerable
 - F, A and I: ~70%
 - R: 35%
 - It is the most needed advancement

Results	F	А	Ι	R	Total
Answerable (%)	76%	70%	71%	35%	64%

 Non-answerable questions require a specialist and improvements in the PubChem's documentation

Applying the Questionnaire

Findability	Accessibility	Interoperability	Reusability
 Plenty of information Identifiers Metadata schemas Data and metadata linkage 	 Well-characterized Protocol Security mechanisms 	 Provides information Knowledge representation (Meta)data formats, e.g., XML, TTL, Json Structured vocabularies 	 Provides information Usage license
 Except, Identifier persistence (Meta)data indexed in a searchable resource 	 Lack of information Security information to access the data manually or by a computer agent (Meta)data longevity plan 	 Lack of information Use of FAIR vocabularies Qualified references used to link data and metadata, and vice-versa 	 No information Relevance and accuracy of (meta)data attributes Provenance schemas and attributes Used domain-relevant community standards

Remarks

FAIR characterization of digital objects' properties is the starting point to understand how close they are to the FAIR principles

- -We proposed an improved questionnaire that can be improved
 - Validate the responses with domain experts
 - Validate the questionnaire with FAIR experts
 - Apply the questionnaire in other scenarios

Questionnaires	R1	R2	R3	R4	R5	R6	R7	R8	R9
FAIR Data Self Assessment Tool	0	0			0		Ο	0	Ο
FAIRDat		Ο		Ο			Ο	0	0
FAIRDataBR	Ο	Ο	Ο		0		Ο	Ο	0
Data Sterwardship Wizard					0	0	Ο	Ο	0
FAIR Implementation Profile (FIP)			Ο	Ο	0	0	0	Ο	Ο
Improved Questionnaire									

-Current approaches do not suffice

IBM Research / Oct. 2024 / © 2024 IBM Corporation

Remarks

- -Questionnaires' strengths
 - Essential for
 - -Overall understanding and
 - -Appreciation of the research life cycle

- -Questionnaires' weaknesses
 - Time consuming
 - Requires experience and technical skills
 - Carries difficulties when inspections is needed
 - Does not scale for several digital objects

Automated Tools for FAIRness Assessment

- -Strengths
 - Performs evaluation without human intervention
 - Scale when evaluating several digital objects
 - More objective
 - Allow comparison of distinct digital objects

- -Weaknesses
 - Requires precise definition of metrics and evaluation tests
 - May be difficult to fit if community standards are not defined
 - May result on using domain-agnostics concepts
 - May not fit community needs

IBM Research / Oct. 2024 / © 2024 IBM Corporation

*Azevedo, L. G., Banaggia G., Tesolin J., Cerqueira R. "An Appraisal of Automated Tools for FAIRness Evaluation". In: 4th Workshop on Metadata and Research 23 (objects) Management for Linked Open Science (DaMaLOS 2024), 2024. DOI: https://repository.publisso.de/resource/frl:6483276

Automated Tools for FAIRness Assessment

Analysis of automated tools for FAIRness assessment

- -Search for existing tools in the literature
 - Discover the tools
 - Elicit requiments
- Examine tools regarding elicited requirements

Literature Review

Abbreviated systematic literature review

Research questions

- RQ1: What are the existing automated tools for FAIRness evaluation?
- RQ2. Which requirements do these tools meet?
- Search string ("Tool" OR "Automated") AND ("Assessment" OR "Evaluation") AND ("FAIRness" OR "FAIRification") AND ("FAIR Principles" OR "FAIR Data")

Search on Scopus, IEEE and ACM digital libraries

-32 works found

- Exclution and inclusion criteria endup with
 - Krans et al. (2022)
 - Peters-Von Gehlen et al. (2022)
 - Slamkov et al. (2022)
 - Sun et al. (2022)
- Gaps on exiting works
 - Abstract characterization and comparison of tools
 - Do not propose or use requirements

Literature Review: Tools

RQ1: What are the existing automated tools for FAIRness evaluation?

- -Tools referenced in the works
 - Krans et al. (2022)
 - Peters-Von Gehlen et al. (2022)
 - Slamkov et al. (2022)
 - Sun et al. (2022)

Search for existing tools in the literature

Tool	Automated?
F-UJI	Yes
FAIR Evaluator	Yes
FAIR Enough [*]	Yes
FAIR-Checker	Yes
ARDC's FAIR Data Self Assessment Tool	No
Checklist for Evaluation of Dataset Fitness for Use	No
CSIRO's 5°Oz Data tool	No
DANS's SATIFYD	No
Data Stewardship Wizard	No
EUDAT's Checklist	No
FAIRdat	No
FAIRenough	No
FAIRshake	No
GARDIAN	No
RDA's Simple Grid	No
Semi-automated workflow for FAIR maturity indicators	No

IBM Research / Oct. 2024 / © 2024 IBM Corporation

* FAIR Enough was found when we looked for a reference of FAIRenough. FAIR Enough is an automated tool based on F-UJI and FAIR Evaluator.

Literature Review: Requirements

RQ2. Which requirements do these tools meet?

- -Requirements
 - Guide the appraisal and development of tools
 - Crucial for making objective FAIRness evaluations and improving digital objects
- Requirements elicited from
 - The works (Krans et al., Peters-Von Gehlen et al., Slamkov et al., and Sun et al.)
 - Tools documentation (F-UJI, FAIR Evaluator, FAIR Enough, FAIR Checker)

Elicited requirements (23 requirements)

Req	Requirement: The tool should			
R1	be fully automated.			
R2	give a FAIRness score /grade.			
R10	be customizable according to the type of digital object and community.			
R12	provide a visual representation (e.g., a badge) of the FAIR assessment results.			
R14	rely on FAIR-enabling services.			
R15	offer guidance on how it is used (e.g., providing user manual, help, and publications).			
R18	disclose its rating system (e.g., evidences and rationale).			
R19	be informative, i.e., teach the user about FAIR.			
R20	give recommendations on how to improve the FAIRness of the evaluated resource.			
R23	support versioning of FAIRness assessment.			

Appraisal of the Tools

Evaluation by reading tools' documentation

-Web pages

–GitHub pages

-Papers

Examine tools regarding elicited requirements

Req	Keyword	F-UJI	FAIR Evaluator	FAIR enough	FAIR Checker
R1	Automated				
R2	Score				
R10	Customizable	\bigcirc	\bigcirc	\bigcirc	
R12	Badge		0	0	
R14	FAIR- enabling services	•	•	•	-
R15	Guidance			0	
R18	Rating system	•		•	
R19	Teach			\bigcirc	\bigcirc
R20	Recommenda tions	Ο	0	0	
R23	Versioning	0	0	0	0
	Labels: Requirer Requirer Requirer	nent is totall nent is not s nent is partia	y supported upported ally supported		

IBM Research / Oct. 2024 / © 2024 IBM Corporation

FAIRness Score



R1. The tool should give a FAIRness **score**/grad.

Without a numeric score, it is difficult to objectively compare results of FAIRness evaluation executed in the same contexts (e.g., FAIR aspects considered, configurations of metrics and tests used in the evaluation).

Main question: How to compute the FAIRness grade for dimensions, principles, metrics, and tests in a way the grade is not only a number but considers aspects like priorities and weights defined by a community and in transparent way for the user?

IBM Research / Oct. 2024 / © 2024 IBM Corporation

FAIRness Badge

R12. The tool should provide a visual representation (e.g., a **badge**) of the FAIR assessment results.



Without a badge, the user does not have the whole assessment in a visual representation.

!? Main question: What is the best representation that present the results' overview for all evaluation levels (principles, metrics, tests)?

Recommendations

R20: The tool should give recommendations on how to improve the FAIRness of the evaluated resource



Without giving recommendations (e.g., recipes or standard schemas), one misses the opportunity to increase the FAIRness of data

- **X** F-UJI, FAIR Evaluator, FAIR Enough: present a log of the execution without explicit recommendations
- **V** FAIR Checker: a set of recommendations for FAIRness improvements with links to training resources, such as FAIR-Cookbook
- **!?** Main question: How to present FAIRness improvements recommendations to be followed by the non-technical users?

Customization

R10: The tool should be customizable according to the type of digital object and community



F-UJI and FAIR-Checker

- Do not support user friendly configuration
- Require software development skills to develop and add new tests in the tool



- 👍 Allow users to group tests in a collection
- Require software development skills to develop and add new tests in the tool

Without the ability to customize the tool, evaluation is limited to agnostic parameters, i.e., does not handle community-specific needs.

!? Main question: How to create FAIRness assessment tools that is easily adaptable by non software development users?

Automated tools Appraisal Results

Tools analysis

- -Similar responses for 15 requirements
 - R1 to R8, R10, R13, R16, R18, R19, R22, and R23
- -Different responses for 8 requirements
 - R9, R11, R12, R14, R15, R17, R20, R21

Fulfillment

- -74%: F-UJI
- -70%: FAIR Checker
- -61%: FAIR Evaluator and FAIR Enough

Tools main strenghts

- Employ good software development practices
- -Use state-of-the art technologies in
 - Software Engineering
 - Semantic Web

Tools main weaknesses

- Reporting features should be improved
- Storage of results and versioning are not implemented

Remarks about the tools

(questionnaires and automated tools)

Requirements are a base for appraising tools

No tool meets all the requirements and stands out as state-of-the-art

Req	An automated tool should						
R1	be fully automated.						
R2	give a FAIRness score/grade.	Req	Keyword	F-	FAIR	FAIR	FAIR
R10	be customizable according to the type of digital object and community.	R1	Automated			enough	
R12	provide a visual representation (e.g., a	R2	Score			$\overline{}$	
	badge) of the FAIR assessment results.	R10	Customizable	•		•	
R14	rely on FAIR-enabling services.	R12	Badge		0	0	
R15	offer guidance on how it is used (e.g., providing user manual, help, and publications).	R14	FAIR- enabling services	•	•	•	•
R18	disclose its rating system (e.g., evidences and rationale).	R15	Guidance	•	•	0	•
R19	be informative i.e. teach the user about	R18	Rating system				-
NT)	FAIR.	R19	Teach	0		-	-
R20	give recommendations on how to	R20	Recommenda tions	0	0	0	•
	resource.	R23	Versioning	0	0	0	0
R23	support versioning of FAIRness assessment						

Req	A questionnaire tool should
R1	have open-ended questions
R2	have questions for data and metadata
R3	provide examples that help answer it
R4	deal with all the FAIR principles
R5	allow for a customization
R6	be intended to generate a FAIRness grade
R7	meet indicators or metrics like RDA FAIR Data Maturity Model
R8	support FAIRness test automation
R9	elicit <mark>evidence</mark> that back the assessment



IBM Research / Oct. 2024 / © 2024 IBM Corporation

Remarks about the tools

(questionnaires and automated tools)

No tool meets all the requirements and stands out as state-of-the-art

-Choosing the best tool is challenging

- There is room to solve the gaps by
 - Evolving existing tools
- or
 - Developing a new tool

To make a choice of a tool

- -Start by
 - Using requirements, like the ones we proposed
 - Identifying the most critial needs
 - Reading the details of our appraisals
- -Then
 - Understand the difficulties to customize an existing tool
 - Test the tools in practice
- Make a decision
 - To use or improve a tool or develop your own



Tools to characterize FAIR and measure FAIRness

Questions? Suggestions?

Leonardo Guerreiro Azevedo IBM Research – Brazil Iga@br.ibm.com



