



**Oberlinnovation** 

How to assess FAIRness of your data – a summary of testing two FAIR validators

Michael Rusongoza Muzoora

Research Associate *Berlin Institute of Health at Charité* 







# Agenda

- FAIR Health Data
- Methods
- Results
- Discussion
- Conclusion







### Disclosure

- We do not have a conflict of interest with any of the mentioned services provided by <u>F-UJI Tool</u> and <u>FAIR-</u>
   <u>Checker</u>
- F-UJI Tool and FAIR-Checker were contacted and are aware of our paper



## FAIR Health Data

- Quality data (exams, laboratory tests, research) needed for prevention, diagnosis and treatment of diseases
- To extract value from data, it must be:
- Why is FAIR data important?
  - Objective of FAIR healthcare data (European Commission)
  - Exchange of FAIR data within European Health Data Space (EHDS)
- How FAIR is my data?
  - Need for FAIR data validation tools



# NFO23

8 – 12 JULY 2023 | SYDNEY, AUSTRALIA

### **FAIR Principles**

#### Box 2 | The FAIR Guiding Principles

To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
- A1.1 the protocol is open, free, and universally implementable
- A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- 12. (meta)data use vocabularies that follow FAIR principles
- 13. (meta)data include qualified references to other (meta)data

#### To be Reusable:

- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
- R1.1. (meta)data are released with a clear and accessible data usage license
- R1.2. (meta)data are associated with detailed provenance
- R1.3. (meta)data meet domain-relevant community standards

Source: Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A et al. Comment: The FAIR Guiding Principles for scientific data management and stewardship. Scientific data. 2016 Mar 15:3. 160018. https://doi.org/10.1038/sdata.2016.18





### F-UJI – An Automated FAIR Data Assessment

- Developed by the FAIRsFAIR project, launched in 03/2019
- FelWebeba</u>sed RESTful service; FAIRness assessment based on metrics (17)
- Source code available on GitHub (free license)

Metrics	
FsF-F1-01D	Data is assigned a globally unique identifier
FsF-F1-02D	Data is assigned a persistent identifier
FsF-F2-01M	Metadata includes descriptive core elements to support data findability.
FsF-F3-01M	Metadata includes the identifier of the data it describes.
FsF-F4-01M	Metadata is offered in such a way that it can be retrieved by machines.
FsF-A1-01M	Metadata contains access level and access conditions of the data.
FsF-A1-02M	Metadata is accessible through a standardized communication protocol
FsF-A1-03D	Data is accessible through a standardized communication protocol
FsF-A2-01M	Metadata remains available, even if the data is no longer available.
FsF-I1-01M	Metadata is represented using a formal knowledge representation language.
FsF-I1-02M	Metadata uses semantic resources.
FsF-I3-01M	Metadata includes links between the data and its related entities.
FsF-R1-01MD	Metadata specifies the content of the data.
FsF-R1.1-01M	Metadata includes license information under which data can be reused.
FsF-R1.2-01M	Metadata includes provenance information about data creation or generation.
FsF-R1.3-01M	Metadata follows a standard recommended by the target research community of the data.
FsF-R1.3-02D	Data is available in a file format recommended by the target research community.
Madified from Lune //	

Modified from: https://www.f-uji.net/index.php?action=methods



## FAIR-Checker

- Developed by the Interoperability working group at French Institute of Bioinformatics
- Web-based tool; 4-step process
  - Semantic web technologies used to check for use of standards

### FAIR-Checker

- 1. Knowledge graph build from semantic annotations (in metadata)
- 2. Existent knowledge graphs used to complete graph
- 3. Test if classes/properties in graph are recognized using OLS/LOV/Bioportal
- 4. Validation of metadata against Bioschemas community profiles







## Methods

1. Identify available **FAIR data validation tools** through literature search

FAIRassist website

2. Select two validation tools for **testing** 

F-UJI – Automated FAIR Data Assessment Tool & FAIR-Checker

3. Test both tools **using three variant data files** with different formats (JSON, TXT, CSV)





### Methods – Resources tested (1/3)

Example FHIR® Observation resource of a detected genetic variant



#MEDINF023

- Fast Interoperability Healthcare Resources (FHIRR)
- Standard exchange format (JSON/XML) for healthcare data
- Structures information into building blocks = "resources"

Source: German Medical Informatics Initiative (MII)

Format: JSDN





### Methods – Resources tested (2/3)

2. Example of a Single Nucleotide Polymorphism (SNP) in the BRCA gene

GDC_Aliquot	Chromosome	Start	End	Num_Pro	bes	Segment	Mean		
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	62920	3298168	737	-0.0858		
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	3301765	1554782	9	7446	-0.1723	
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	15551554	4	16544783	1	452	-0.1094
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1656913	2	16583114	Ļ	3	-1.7836
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1660483	4	16782786	;	46	-0.364
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1678480	1	16788139	)	5	-1.7637
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1685594	2	17010382	2	82	0.0859
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1701438	3	17016429	)	4	-1.6115
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	1701858	8	21992508	1	3311	-0.1641
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	2199666	4	22001786	;	2	-2.9341
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	2200202	5	25230200	)	1907	-0.1583
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	2523206	9	25320455	;	35	0.3834
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	2533681	9	28740003	\$	1413	-0.1428
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	2875315	9	34626073	1	3532	0.2274
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	3463705	3	34649941		22	-0.41
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	3465096	9	61661425	;	16476	0.213
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	6166191	2	61676223	1	15	0.6696
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	6167642	5	65419674	ļ.	2760	0.2119
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	6541988	6	65442811		14	0.6801
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	6544444	2	72284676	)	4480	0.208
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	7229043	9	72297687	,	9	-0.6471
e50f848e-a469	-41c1-b44b-4dbc28	3fd0a6	1	7230273	5	72345465		48	1.0941



• SNP: genomic variant at a single base position in DNA
• <u>Source: National Cancer Institute, Genomic Data Commons (GDC) Data</u> <u>Portal</u>
<ul> <li>Project: TCGA-BRCA, Breast Invasive Carcinoma</li> </ul>
• <u>Format: TXT</u>





### Methods – Resources tested (3/3)

### 3. Breast Cancer Gene Expression Profiles (METABRIC) dataset

∞ patient_id =	# age_at_diagnosis \Xi	▲ type_of_breast_s =	▲ cancer_type =	▲ cancer_type_deta =	▲ cellularity =	# chemotherapy
0 729	21.9 96.3	MASTECTOMY 59% BREAST CONSER 40% Other (22) 1%	Breast Cancer 100% Breast Sarcoma 0%	Breast Invasive Du 79% Breast Mixed Duct 11% Other (197) 10%	High         49%           Moderate         37%           Other (254)         13%	0
0	75.65	MASTECTOMY	Breast Cancer	Breast Invasive Ductal Carcinoma		0
2	43.19	BREAST CONSERVING	Breast Cancer	Breast Invasive Ductal Carcinoma	High	0
5	48.87	MASTECTOMY	Breast Cancer	Breast Invasive Ductal Carcinoma	High	1
б	47.68	MASTECTOMY	Breast Cancer	Breast Mixed Ductal and Lobular Carcinoma	Moderate	1
8	76.97	MASTECTOMY	Breast Cancer	Breast Mixed Ductal and Lobular Carcinoma	High	1
10	78.77	MASTECTOMY	Breast Cancer	Breast Invasive Ductal Carcinoma	Moderate	0
	56.45	BREAST CONSERVING	Breast Cancer	Breast Invasive Ductal Carcinoma	Moderate	1

• Source: Molecular Taxonomy of Breast

Cancer International Consortium

### (METABRIC)

- <u>Dataset</u>: by Prof. C. Caldas, from <u>BioPortal</u>
- <u>Format</u>: CSV







### **Result 1:** Example FHIR® *Observation* resource, MII (JSON)





@TheInstituteDH



### **Result 2:** Example SNP, NCI (TXT)





@TheInstituteDH





### **Result 3:** Gene expression profiles, Kaggle (CSV)





@TheInstituteDH



# Summary

### F-UJI Tool

<u>JSON file</u>: FAIRness low (14%); contrary to expectations <u>TXT file</u>: FAIRness low (4%); lower than expected <u>CSV file</u>: moderately FAIR (56%)

- ightarrowTool detected characteristics specific to all 4 FAIR principles in files
- $\rightarrow$ Current version not suitable for FHIR JSON format







# Summary

### FAIR-Checker

<u>JSON file</u>: Unique IDs (F1A) & open resolution protocol (A1.1) detected, others failed <u>TXT file</u>: Unique IDs (F1A) & open resolution protocol (A1.1) detected, others failed <u>CSV file</u>: Meets findability, accessibility and interoperability principles, but 2/3 reusability principles

 $\rightarrow$ Tests highlight reusability & interoperability of data can be ensured through:

- Use of RDF-compliant metadata &
- Providing license and provenance information within metadata





## Discussion

- Tests performed with F-UJI and FAIR-Checker are only indicative of the true potential of the tools; they are still in development (demo versions were used)
- 2. It would be useful if **guidance on expected resource format** were provided to users
- Facilitating FAIRness of genomic data could be assisted if tools were compatible with genomic file formats (e.g., VCF)
- Output of the tools differed and was dependent on the submitted file format.
- Only one of the tools provided an aggregate FAIR score, both tools recorded scores for each FAIR principle



### References

- M. D. Wilkinson et al., 'The FAIR Guiding Principles for scientific data management and stewardship.', Sci. Data, vol. 3, p. 160018, Mar. 2016, doi: 10.1038/sdata.2016.18.
- European Commission and Directorate-General for Research and Innovation, Turning FAIR into reality : final report and action plan from the European Commission expert group on FAIR data. Publications Office, 2018. doi: 10.2777/1524.
- FAIRassist and University of Oxford, 'FAIRassist.org', Jul. 02, 2019. https://fairassist.org/#!/ (accessed Nov. 14, 2022).
- Medizininformatik Initiative, 'Genetic variant assessment', 2022. https://simplifier.net/medizininformatikinitiative-modulomics/example-miimolgen-variante-1 (accessed Nov. 26, 2022).
- National Cancer Institute, 'GODFS\_p\_TCGA\_bII7\_II8\_SNP\_N\_GenomeWideSNP\_6\_ED6\_778DI4.grch38.seg.v2.txt', Aug. 23, 2018. (accessed Nov. 26, 2022).
- 'Breast Cancer Gene Expression Profiles (METABRIC)', kaggle, 2019. https://www.kaggle.com/datasets/raghadalharbi/breast-cancer-gene-expressionprofiles-metabric (accessed Nov. 26, 2022).
- Anusuriya Devaraju and Robert Huber, 'F-UJI An Automated FAIR Data Assessment Tool (v1.0.0)', Zenodo, 2020. doi: 10.5281/zenodo.4063720
- Rosnet, Thomas, Lefort, Vincent, Devignes, Marie-Dominique, and Gaignard, Alban, 'FAIR-Checker, a web tool to support the findability and reusability of digital life science resources', Jul. 2021, doi: 10.5281/ZEN0D0.5914307.
- Robert Huber et al., 'FAIRsFAIR Data Object Assessment Metrics\_v0.4\_PublicFeedback', Google Docs, Oct. 12, 2020.
   <a href="https://docs.google.com/document/d/lymkzVmF\_BJmKT0200SR01Y0JaPxefJJ284AKUJUIGeM/edit?usp=sharing&usp=embed\_facebook">https://docs.google.com/document/d/lymkzVmF\_BJmKT0200SR01Y0JaPxefJJ284AKUJUIGeM/edit?usp=sharing&usp=embed\_facebook</a> (accessed Nov. 17, 2022).
- FAIRsFAIR "Fostering FAIR Data Practices In Europe", 'FAIRsFAIR Data Object Assessment Metrics: Request for comments', FAIRsFAIR, Jul. 10, 2020. <u>https://www.fairsfair.eu/fairsfair-data-object-assessment-metrics-requestcomments</u> (accessed Nov. 17, 2022).
- J. D. B. Jacobsen et al., 'The GA4GH Phenopacket schema defines a computable representation of clinical data', Nat. Biotechnol., vol. 40, no. 6, Art. no. 6, Jun. 2022, doi: 10.1038/s41587-022-01357-4.





### Thank You!



### Michael Rusongoza Muzoora *michael.muzoora®bih-charite.de*



Caroline Stellmach caroline.stellmach@bih-charite.de

