

1 Testing the FAIR metrics on data catalogs

2 *“Metrics, Metrics can you recall,*

3 *Which data catalog is the FAIRest of them all?”*

4

5 [Jarno A A van Erp](#)¹, [Carolyn D Langen](#)¹, [Anca Boon](#)¹, [Kees van Bochove](#)¹

6 ¹The Hyve B.V., Arthur van Schendelstraat 650, 3511 MJ Utrecht

7

8 *Corresponding Author:*

9 Kees van Bochove¹

10 Arthur van Schendelstraat 650, 3511 MJ Utrecht

11 Email address: office@thehyve.nl

12

13 **Abstract**

14 The introduction of the FAIR –Findable, Accessible, Interoperable, Reusable– principles has
15 caused quite an uproar within the scientific community. Principles which, if everyone adheres to
16 them, could result in new, revolutionary ways of performing research and fulfill the promise of
17 open science.

18 However, to bring about these changes, data users need to rethink the way they treat scientific
19 data. Just passing a dataset along, without extensive metadata will not suffice anymore. Such
20 new ways of executing research require a significantly different approach from the entire
21 scientific community or, for that matter, anyone who wants to reap the benefits from going FAIR.

22 Yet, how do you initiate this behavioral change? One important solution is by changing the
23 software scientists use and requiring data owners, or data stewards, to FAIRify their dataset.

24 Data catalogs are a great starting point for FAIRifying data as the software already intends to
25 make data Findable and Accessible, while the metadata is Interoperable and relying on users to
26 provide sufficient metadata to ensure Reusability. In this paper we analyse to what extent the
27 FAIR principles are implemented in several data catalogs.

28 To determine how ‘FAIR’ a dataset is, the FAIR metrics were created by the GO-FAIR initiative.
29 These metrics help determine to what extend data can be considered FAIR. However, the metrics
30 were only recently developed, being first released at the end of 2017.

31 The Hyve has tested/evaluated three popular open source data catalogs based on the FAIR
32 metrics: CKAN, Dataverse, and Invenio. Most data stewards will be familiar with at least one of
33 these.

34 Within this white paper we provide answers to the following questions:

- 35 • Which of the three data catalogs performs best in making data FAIR?
- 36 • Which data catalog utilizes FAIR datasets the most?
- 37 • Which one creates the most FAIR metadata?
- 38 • Which catalog has the highest potential to increase its FAIRness, and how?
- 39 • Which data catalog facilitates the FAIRification process the best?

40 Introduction

41 Earlier this year the international organisation GO-FAIR opened its first offices in Leiden, the
42 Netherlands and Hamburg, Germany. The organisations aim is to promote and facilitate the
43 implementation of the “internet of FAIR data and services”. Several countries, including
44 Germany and the Netherlands, have committed themselves to implementing an infrastructure
45 capable of supporting the newly arisen needs introduced by the FAIR data principles. This
46 includes clearly defined data permissions (to ensure the research data is not misused), addition of
47 extensive metadata, noninvasive data sharing, and eventually transmission and integration of
48 data/information across different organisations in different countries with different laws and
49 regulations. The efforts of GO-FAIR should result in a clear understanding across organisations
50 what the requirements for each dataset are. Funding organisations, such as NWO and H2020,
51 have already integrated requirements for FAIR data in their research grant application process, in
52 this way safeguarding that data created with their funds will be reusable for future research.

53 One major concern regarding these new requirements is: How to facilitate making research data
54 FAIR? FAIRifying data could become a time consuming activity when the FAIR data stewards
55 have not sufficient insight in the research question or needs to revise a large amount of data. This
56 increases the need for tools to identify FAIR business practices which require a low effort input
57 while resulting in high value output. The tools that are most suitable for these low effort/high
58 value FAIR business practices are data catalogs, as they are already used to make data findable
59 and accessible, ensuring metadata is interoperable, and facilitating reuse of data.

60 Data catalogs

61 Combined with the vast amounts of data that need to be processed in scientific and medical
62 research nowadays, one important system requirement is that data does not get lost. A common
63 approach for preventing data loss is storing it in a data catalog. Data catalogs help to organize,
64 structure and track metadata and data generated, so that the information can be saved and shared
65 within an organisation. The use of data catalogs could even result in scientists getting more
66 citations, as they create opportunities to elaborate on or reuse prior research. For example, a data
67 catalog makes it much easier to search for relevant data.

68 To enable the entire scientific community to fully benefit from research data, the reusability of
69 data should be improved in a trustworthy manner, protecting both the data producer and the

70 external data re-user. By improving the quality and comparability of research data, fellow
71 scientists should be able to reuse a particular dataset. Establishing trust between data producers
72 and external data re-users is an issue requiring a stronger behavioural change among scientists
73 than just asking them to add extra metadata to datasets. To facilitate the necessary change,
74 several funding agencies nowadays require that grant recipients provide a Data Management
75 Plan or Data Stewardship Plan, which includes a description of how the data will be made
76 available to fellow researchers....

77 Important aspects of data management with regard to the reusability of data are the location and
78 method by which the data is stored. Besides, the risk of someone misusing sensitive data should
79 be minimized and this aspect should be duly considered when making data available publicly or
80 within an institute. Luckily, there are software solutions available that limit the risks of misuse,
81 while ensuring the data can be reused by the original creator, the institute, and potentially anyone
82 around the world. Ultimately, data reuse increases the value of datasets thus increasing the
83 likelihood that (public) money is being spent in an efficient manner.

84 FAIR principles

85 Optimal data reusability is at the core of the FAIR principles. It ensures that humanity will get
86 the most value out of publicly funded scientific research. In order to achieve this, the FAIR
87 principles require metadata additions, data point annotations, and precise descriptions of
88 potentially relevant information, thereby ensuring findability, accessibility, and safeguarding
89 reusability by others. All these qualities combined make data catalogs the most attractive tools to
90 start FAIRifying data.

91 However, their popularity means scientists can choose between a number of options. Many data
92 catalogs have implemented changes to better adhere to the FAIR principles. It is rarely clearly
93 stated, though, to which metrics they adhere and to what extend.

94 **Background**

95 To come up with a method of determining FAIRness of catalogs, we first need to define what
96 exactly a data catalog is. Subsequently, we will discuss the method used to determine how FAIR
97 the three catalogs are that The Hyve analysed.

98 Data catalogs

99 According to the Data Catalog vocabulary (DCAT), published by the World Wide Web
100 Consortium (W3C), a data catalog is a curated collection of metadata about datasets. Practically
101 speaking, this means a service which allows users to deposit, describe, share, discover and
102 explore datasets. When data is correctly and accurately curated, they can be better understood
103 and (re)used, which increases the value of any given dataset.

104 The Hyve evaluated three data catalogs using the FAIR metrics: CKAN, Dataverse and Invenio.
105 The capabilities we explored include:

- 106 • Recognition of a variety of file formats
- 107 • Digital Object Identifier (DOI) generation
- 108 • Fast search indexing (e.g. using open source search engines SOLR or ElasticSearch)
- 109 • Harvesting (meta)data from external catalogs

110 The unique features of each catalog are described below.

111 CKAN

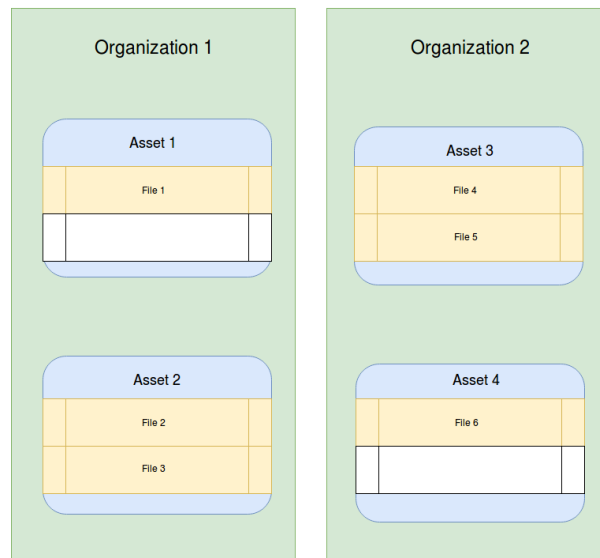
112 The CKAN (Comprehensive Knowledge Archive Network) catalog allows for the creation of
113 “organizations” – entities which supply data. Organizations can contain multiple datasets, called
114 assets within CKAN. An asset combines one or more files with metadata and tags. Views can be
115 attached to the organisations allowing users to preview the assets using maps, graphs and tables
116 (see Figure A). A news feed shows users recently added and/or modified assets.

117 CKAN has a plug-in system for adding features such as enhanced monitoring, custom user roles
118 and permissions. At The Hyve we developed a plug-in expanding these options, focussing on a
119 more fine grained accessibility mechanism within CKAN. The plug-in system adds custom
120 metadata fields to datasets and allows users to search the data files. There is also a plug-in
121 available that allows CKAN to expose and explore metadata from other catalogs using RDF
122 documents serialized with DCAT and export CKAN data as RDF DCAT endpoints. This creates
123 the option to register the catalog as a *FAIR data point*.

124 CKAN stores its metadata in a PostgreSQL database. Files uploaded to CKAN can be stored on a
125 local drive, a network connected drive, or on Cloud storage solutions such as S3, Azure or
126 Google Cloud Storage.

127

CKAN Overview



128

129 **Figure A: CKAN structure.** CKAN allows for multiple organizations. Organizations consist of assets which contain one or more
130 files.

131 Dataverse

132 Dataverse works with data repositories which are called *dataverse*. The catalog allows users to
133 create a dataverse within a dataverse, where each dataverse has its own administration rights. As
134 such, read and write permissions of each dataverse (and its datasets) can be controlled
135 independently, and metadata can be assigned to a dataverse, a dataset, or to a single file. The
136 recursive structure of a dataverse is illustrated in Figure E. Dataverse uses a PostgreSQL
137 database, combined with Solr for searching, and a local file storage system for saving files.

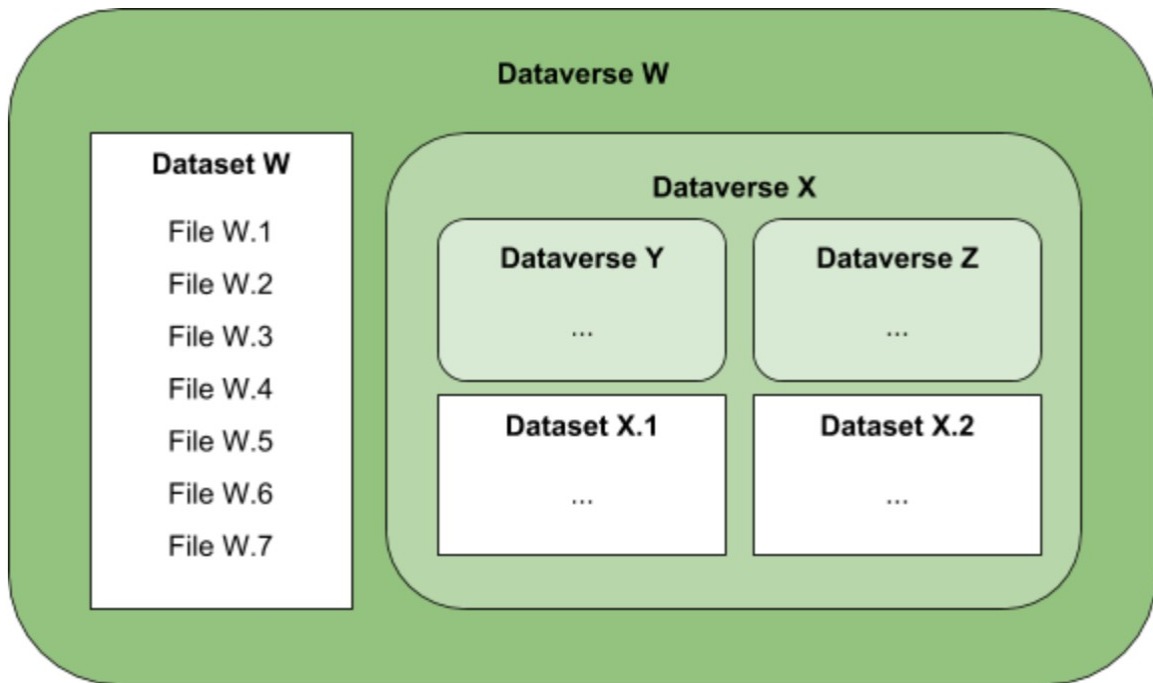
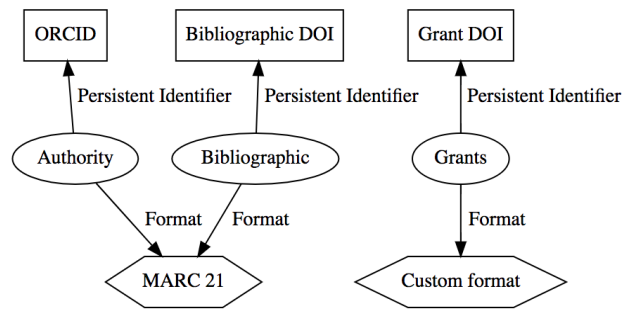


Figure B: Dataverse structure. *Dataverse is a data catalog containing datasets and possibly a number of Dataverses. Each dataset contains files. In this example, Dataverse W contains a dataset with seven files and Dataverse X, consisting of Dataset X.1 and X.2, Dataverse Y and Dataverse Z.*

138 Invenio

139 Invenio is a data catalog developed by CERN, the European Organization for Nuclear Research,
 140 to share their data publicly with fellow scientists. After more than fifteen years of experience
 141 with Invenio, CERN developers released a new, modularly-structured version, with three types
 142 of modules defined as base, core and additional feature. All modules are available in the Python
 143 Package Manager (PyPM) as separate components, which can be replaced by custom-made
 144 solutions.

145 The data model of Invenio (Figure X) consists of linking DOIs with a JSON (JavaScript Object
 146 Notation) Schema representation for the associated metadata. This grants a certain freedom to
 147 create links between datasets while at the same time limiting the complexity to a predefined
 148 model.



149

150 *Figure C: Data model of Invenio. Reproduced from: [http://invenio.readthedocs.io/en/latest/developersguide/create-a-](http://invenio.readthedocs.io/en/latest/developersguide/create-a-datamodel.html)*
 151 *[datamodel.html](http://invenio.readthedocs.io/en/latest/developersguide/create-a-datamodel.html)*

152 What is FAIR?

153 As mentioned before, making data FAIR –Findable, Accessible, Interoperable and Reusable–
 154 requires a behavioral change from scientists. Many researchers recognize that unwillingness to
 155 share results, loss of data and a focus on publication output instead of research quality are
 156 detrimental to the research community. By rating the data quality researchers produce rather than
 157 focussing on the ranking of the journal in which results are published, the GO-FAIR initiative
 158 aims to change the way scientific data is valued and wants to facilitate the process of behavioral
 159 change.

160 Ensuring that data can be reused creates a higher value proposition for generating data. By
 161 describing data with rich metadata and annotating the dataset itself, a computer, and eventually
 162 every machine, will be able to interpret the data: the machine “knows” what information the
 163 dataset contains, can link similar data, is able to create a knowledge graph out of these links. This
 164 would, all in all, reduce the time researchers spent searching for potentially interesting datasets.
 165 Besides, machines will be able to convert any created knowledge graph to human readable
 166 formats, enabling researchers to explore and use these more easily for research purposes. With
 167 the machine being able to identify what is inside any given dataset, it can make the data
 168 interoperable with other datasets. This leads to increased analytical capability and improved data
 169 maintenance.

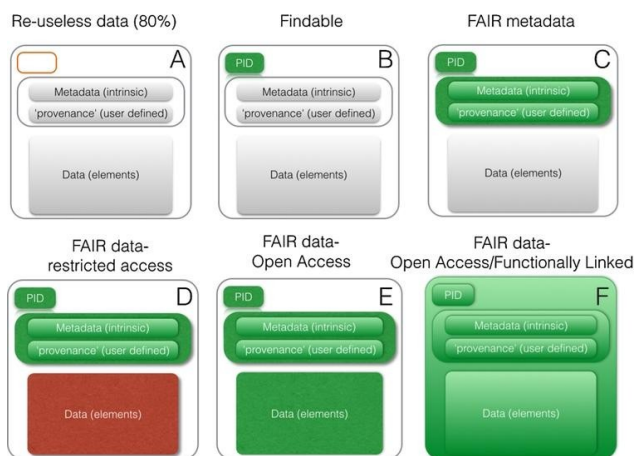
170 To be able to assess the FAIRness of data, the FAIR metrics were developed. The metrics consist
 171 of multiple rules describing what is needed to comply with each one of the FAIR principles. For

172 example, to be fully compliant to the first principle of Findability (F1) a Uniform Resource
 173 Locator (URL) to a registered identifier scheme should be provided, along with a URL that links
 174 to a document containing the relevant policy for when this identifier scheme becomes
 175 deprecated.

176 Degrees of FAIRness

177 The release of the FAIR principles created confusion about how to actually implement these
 178 standards. Degrees of FAIRness were introduced to bring clarity and explain if data needs to
 179 adhere to all criteria of a FAIR metric to even be considered FAIR. In this paper, we make
 180 suggestions how to use these gradients with regards to FAIR and end-user software (data
 181 catalogs in this case), contradicting a common perception that tools are either FAIR or not FAIR.

Data as increasingly FAIR Digital Objects



182

183 A data owner can determine to what extent he wants his data to be FAIR. For example, privacy
 184 sensitive patient data are never meant to be freely accessible. Figure D shows the guidelines that
 185 scientists should adhere to in order to obtain a certain degree of FAIRness in their datasets.

186

187 **Figure D: Increasingly FAIR digital objects.** To be considered FAIR, some steps need to be taken. This figure gives an overview
 188 how different measures increase the FAIRness of a dataset.

189 However, there is some discussion in the field if Figure F is a proper representation and
190 description of the degrees of FAIRness, as a dataset for internal use can comply to only a number
191 of FAIR principles and still be used in a FAIR manner. To make a dataset Findable to a degree, it
192 could be sufficient to add just a PID (Persistent Identifier), metadata in a machine-readable
193 format, and provenance.

194 Despite these shortcomings, the FAIR metrics do help determine the level of FAIRness of a
195 certain dataset or data catalog. To truly be FAIR, all data and metadata should be stored in an
196 online repository. Otherwise, how can someone else find and access them?

197 FAIR metrics

198 The FAIR metrics consist of four separate “groups”, making the standard distinction between
199 Findable, Accessible, Interoperable and Reusable, where each group has different metrics to
200 determine if data adheres to the corresponding principle. This means that FAIR metrics FM-A1.1
201 corresponds to FAIR principle A1.1. This naming convention will be used in this paper, with the
202 addition of FM referring to the metric related to the principle.

203 One aspect that needs to be emphasized, is that the FAIR metrics will be constantly changing and
204 evolving with the introduction of new technologies and standards in years to come. Currently, the
205 FAIR metrics have a strong emphasis on determining the FAIRness of a dataset rather than the
206 FAIRness of software. For example, it is not stated where certain information needs to be
207 located. In general, this additional information should be present within the dataset, although
208 from a software perspective it should be enough to only link to certain, standard, information. If
209 the metrics are used to determine the FAIRness of data contained within software, the way to
210 find that data, or its location, needs to be stated clearly. When you want to assess the FAIRness
211 of software and focus on automatic machine-readability, it currently is often unclear where the
212 machines should search for specific information. Implementing this could lead to easier ways of
213 connecting various FAIR tools.

214 As a tool to automatically determine the FAIRness of a data catalog is currently being developed
215 by GO-FAIR, we performed our analyses by hand.

216 **Data catalog FAIRness review**

217 To determine the FAIRness of the three data catalogs, we looked at two ways data could be
218 handled: manual and automated. Manual is defined as requiring additional effort from the
219 reviewer, whereas automated meant that the data catalog can make machine readable exports
220 from the data without the need for additional input from the user.

221 For both manual and automated scores we defined three outcomes per matrix: present, partial or
222 absent. Present meaning that the catalog fully meets the criteria, partial meaning that it meets
223 only part of the criteria, and absent implying it did not or not sufficiently meet the criteria. For
224 each partial outcome, we identified what element was missing. An overview of missing FAIR
225 elements of the three data catalogs can be found in Appendix A.

226 Leveraging the plugin design of CKAN, the ckanext-dcat plugin was added for automated DCAT
227 exports. For Invenio the Zenodo version was used, as the FAIR metrics dataset was published in
228 this version.

229 The overall outcome of the review is represented in Table A. The differences in findability,
230 accessibility, interoperability and reusability between the three data catalogs are presented in
231 Table Y. See Graph A for a visualization of Table Y.

Metric	Dataverse		CKAN		Invenio		Meaning automated	Legend
	Manual	Automated	Manual	Automated	Manual	Automated		
F1A	2	2	0	0	2	2	Automatically created	2 Present
F1B	0	0	0	0	0	0	New IRI automatically findable when old identifier scheme becomes deprecated	1 Partial
F2	2	1	2	1	2	1	Metadata and data are automatically annotated in a machine-readable format	0 Absent
F3	2	2	2	2	2	2	IRI is automatically added to export of metadata	
F4	2	2	2	2	2	2	When (meta)data is made public it is automatically indexed	
A1.1	1	1	1	1	1	1	Protocol description is automatically added	
A1.2	1	1	1	1	1	1	Machine knows if it does not have access and knows how to gain access and can perform the required action	
A2	0	0	2	2	0	0	Metadata automatically stays when data is deleted	
I1	2	1	2	1	2	1	Data is automatically annotated with a knowledge representation language	
I2	2	1	2	1	2	1	Data is automatically annotated with FAIR vocabularies	
I3	1	1	2	1	0	0	Relationships are automatically added to metadata, possible relationships are automatically discovered	
R1.1	0	0	2	0	2	2	A machine automatically knows what it can and cannot do with the data	
R1.2	2	2	2	0	2	2	Able to automatically create an overview of the provenance that meets the minimal requirements*	
R1.3	0	0	0	0	0	0	Automated request for certification from a recognized body	

*R1.2 minimal requirements:

- Who/what/When produced the data (i.e. for citation)
- Why/How was the data produced (i.e. to understand context and relevance of the data)

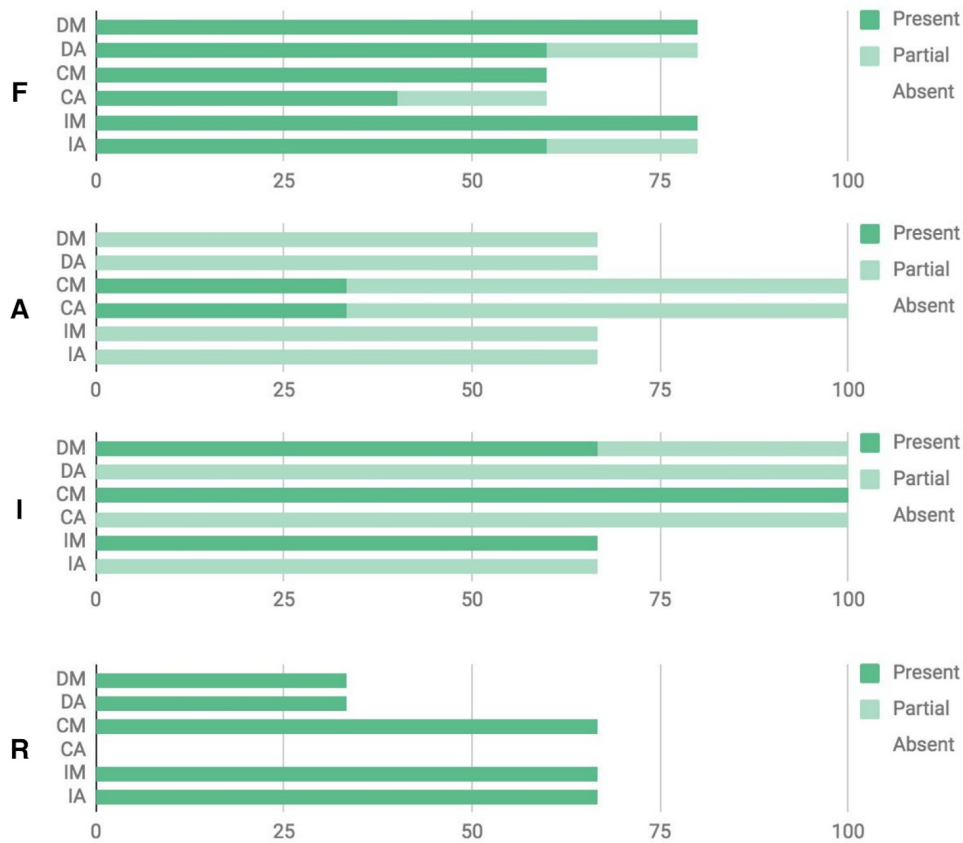
233
234

Table A: Outcome of the FAIR metrics review per metric. The row on the right side specifies what was expected from automated data handling.

235

	Dataverse		CKAN		Invenio	
	Manual	Automated	Manual	Automated	Manual	Automated
Present:	7	4	9	3	8	5
Partial:	3	6	2	6	2	5
Absent:	4	4	3	5	4	4
Total (max=14):	8.5	7.0	10.0	6.0	9.0	7.5

237 **Table B:** General outcome of the FAIR metrics review for data catalogs. A total of 14 points could be granted to each tool. Total
 238 was calculated as follows: $Total = present + (partial/2)$



239

240

Graph G: Visual representation of the FAIR metrics review outcome. See Table A for details.

241

242 Thanks to the liberty CKAN offers, it can be judged the FAIRest of the three catalogs. However,
243 depending on the choices users make, there will be significant differences in the FAIRness of
244 data. A downside of using CKAN (originally an open data repository) is that it offers by default
245 no authentication mechanism, which is an important functionality for scientific communities
246 (especially those dealing with data from human subjects) with laws such as the General Data
247 Protection Regulation (GDPR). Yet despite this shortcoming, CKAN does have the highest score
248 for manual FAIRness.

249 Regarding the FAIRness of Dataverse and Invenio, the major differences are that Invenio
250 provides better permission support. As a downside it does not link datasets adequately. Dataverse
251 includes a license, but it only specifies the licence name. This provides insufficient information.
252 Dataverse does allow the user to create a link between datasets, but this feature can be
253 considered too limited as it provides not enough options to describe meaningful relationships.

254 **Discussion**

255 When we combine the results for both manual and automated generation of FAIR data, CKAN
256 performs worse than both Dataverse and Invenio. Therefore, we conclude that CKAN is good at
257 handling data that is already FAIR. However, the catalog will be less helpful in the data
258 FAIRification process. This draws our attention to different angles that software FAIRness can be
259 viewed from and raises the question if a data catalog is more FAIR when it helps users in making
260 their data FAIR or when it supports data that is already FAIR, without facilitating of the
261 FAIRification process?

262 The FAIR metrics

263 One remarkable outcome of our FAIR review is the low score on the metrics evaluating
264 accessibility, with Dataverse and Invenio both performing worse than CKAN. When looking at
265 the specific criteria the data catalogs did not meet, it becomes clear that CKAN only scores
266 higher because of its ability to preserve metadata when the original dataset has been removed
267 (principle A2). As for the other two metrics, FM-A1.1 and FM-A1.2, all three data catalogs
268 achieve only a partial score. When all data catalogs render the same partial score, it raises the

269 question if this is a result of the FAIR metrics not being defined clear enough or because the
270 proposed methods in the FAIR metrics have not yet been implemented.

271 The question of unclear definition versus incomplete implementation comes to mind especially
272 when assessing the FAIRness of software, as the FAIR principles are not only meant to be
273 guidelines, but also set out a roadmap to what should eventually be implemented. For example,
274 FAIR principle R1.3 states: “(meta)data meet domain-relevant community standards”. The
275 community around data catalogs is the scientific publishing and data preservation community,
276 yet it is currently doubtful if this is the same community that determines what the domain-
277 relevant standards are. Even within a specific field such as the life sciences, there are many
278 standard sets for metadata to choose from (e.g. schema.org, specifications such as DATS and
279 DCAT or the HCLS Community Profile). The associated metric (FM-R1.3) also states that “a
280 valid result” is given, when there is a “successful signature validation”. The result is based upon
281 “a certification saying that the resource is compliant”, which needs to be provided. But at the
282 moment it is unclear who is authorized to give this certification and/or signature. Can someone
283 just make a certificate and say “This is my community, therefore I set the standards and decide
284 that I meet those criteria”? Or is there need for an external body to certify compliance? Should
285 the communities establish the standards and certification bodies themselves with the help of GO-
286 FAIR?

287 Currently, it seems the FAIR metrics outline a prevalent struggle within the FAIR community:
288 inclusiveness or interoperability through standardization? Using standards for interoperability
289 automatically means everyone not using those standards will be excluded. A solution could be to
290 map these standards, an effort already started by FAIRsharing. However, this current mapping
291 project could very well never be finished with new standards emerging that might be more
292 widely adopted. This could be an everlasting discussion. Therefore, we will use what is currently
293 available for reviewing the FAIRness of software and leave this discussion to the FAIR
294 community.

295 Data catalogs and FAIRness

296 Findability

297 The main function of data catalogs is providing an overview of the available data. It is therefore
298 not surprising that all data catalogs score high on the Findability metric. What they all are

299 lacking, though, is a procedure for when the original identifier scheme becomes deprecated. The
300 catalogs also lost half a point for “the availability of machine-readable metadata that describe a
301 digital resource” with regard to automatic creation of FAIR data in certain formats. At this
302 moment, metadata still needs to be added manually by the user.

303 Accessibility

304 The Accessibility metrics has a surprisingly low score. This can be a result of the way the metrics
305 have been formulated or because reality has not yet caught up with the FAIR metrics, depending
306 on which metric is reviewed. Accessibility, for example, has three subcategories: FM-A1.1, FM-
307 A1.2 and FM-A2. While the latter is clear and makes sense: “The existence of metadata even in
308 the absence/removal of data”, the right way to implement the first one is unclear: “The nature
309 and use limitations of the access protocol.”

310 If the secure communication protocol HTTPS is used for data transfer, the metadata export
311 should include a description of the HTTPS protocol as per FM-A1.1. However, if the machine is
312 already familiar with the protocol, does a description still need to be included? FM-A1.2,
313 “Specification of a protocol to access restricted content”, is an example of a metric that has not
314 yet been implemented in any data catalog. Implementation of this metric would require a
315 machine that automatically knows how to access a data source and execute the task needed to
316 gain access.

317 Interoperability

318 For interoperability, all data catalogs reach partial to high scores, depending on whether the
319 automatic or manual score is consulted. This outcome is as expected. Fully automated
320 interoperability would mean automatic detection of what the dataset contains and linking this to
321 an already existing knowledge base. This information would then be used to automatically create
322 links to vocabularies and ontologies and search for relations between datasets. To be able to
323 perform such a complicated task, which needs to be accurate as well, some serious AI power is
324 needed. This, in it’s turn, requires a large number of FAIR datasets to train the AI. Unfortunately,
325 neither of these two requirements are currently met. Therefore, asking the question “How do data
326 catalogs perform on automated interoperability?” is at the moment nonsensical. For now, a better
327 question would be “Does this data catalog link to a machine-readable version upon addition of a
328 dataset?”, when interoperability is a necessity or high priority.

329 Reusability

330 When it comes to reusability the data catalogs all achieve a low score. This finding is disturbing,
331 as one purpose of data catalogs is to ensure the reuse of data. As discussed above, FM-R1.3 is
332 formulated in such a way that full reusability cannot currently be achieved. The other two
333 reusability metrics emphasize licensing and provenance. Although licenses can be added to all
334 three catalogs, within CKAN adding a license is optional. The same goes for provenance. In
335 CKAN the user has the option to specify information. Since this should be a requirement for
336 FAIR data, CKAN achieves a lower score for automated data FAIRification.

337 Manual versus automated

338 All three evaluated catalogs are at least in part capable of handling user-created FAIR data.
339 Manually making a dataset adhere to the the highest FAIRness standards requires a significant
340 effort. A data catalog which supports FAIR data without automation is therefore not a serious
341 contender when deciding which data catalog best meets the user's needs. This is an area where
342 these catalogs can improve upon significantly.

343 Another factor to consider when determining which data catalog to use, is to what extent the tool
344 helps the user to make data FAIR. One of the easiest ways to achieve this is by adding
345 configurable fields, allowing an administrator to decide which fields are mandatory and which
346 are optional. This ensures that the necessary metadata are added, allowing the tool to
347 automatically convert the information into a machine-readable format.

348 Conclusion

349 Although CKAN with manual FAIRification of data by the user has the highest single FAIR
350 metrics score, it does not score highest overall. This is because the catalog does not help and
351 guide users to make data FAIR as much as the other tools. A catalog should help the user with
352 making their data FAIR by default as much as possible.

353 This is where the historical differences between Invenio and Dataverse become apparent.
354 Dataverse was created as a data catalog for researchers, whereas Invenio was developed for
355 storage of bulk data. This results in Dataverse focusing on storing data in such a way that it can
356 be used for publication, whereas Invenio provides a higher quality implementation of the FAIR
357 data principles that ensure trustworthy reusability.

358 Taking everything into account, our overall conclusion is that CKAN can handle data that is
359 already made FAIR better. While Dataverse can be considered just as FAIR as Invenio. However,
360 the ultimate goal of the FAIR initiative is, as mentioned, to change the behavior of researchers
361 and data stewards handling data and have them reconsider how to publish the data they create:
362 ensuring high quality metadata are added and establishing trust regarding reusability by clearly
363 defining rules and developing guidelines for access and permissions. With respect to these
364 aspects, Dataverse is the one that outshines the other two catalogs.

365 Next steps?

366 The low score for Accessibility surprised us. However, bear in mind that score resulted from a
367 number of FAIR requirements that cannot be fully met at the moment. What can we learn from
368 this? In the first place, the low score indicates that the techniques used for granting and retracting
369 access have not yet been fully developed, implemented and accepted for usage according to the
370 FAIR data principles. This does not mean that this aspect is overlooked. As a software developer,
371 you for example may assume that “machine readable accessibility” means that servers and/or
372 clients are automatically authorized and authenticated based upon the identity provided, together
373 with the request made.

374 The high overall scores for Findability reassures that these principles are already widely accepted
375 and implemented. This means that the scientific community is ready for the next step:
376 accessibility. Ideally, this would mean creating an infrastructure where (external) researchers can
377 request data and only need to accompany this request with a link to a verified online identity. A
378 number of scientists have already expressed an interest in supporting external access and identity
379 providers, with techniques such as OAuth and websites like ORCID.

380 Manual vs Automated or Reality vs Future?

381 In this paper, two types of FAIRness were discussed with regards to software: manual and
382 automated. However, this raises the question if it is even realistic to demand fully automated
383 FAIR data creation. Of course, such techniques are not available yet and might not be for some
384 time to come.

385 Therefore, a more sensible question would be: Do we expect researchers, for the time being, to
386 manually add all the metadata needed to make their data FAIR? Can this be a task be left to a
387 machine? Imagine AI being able to define and add FAIR metadata based on headers, column

388 names, data published in research papers, drafts of scientific papers, et cetera. This might not
389 result in full FAIRness but it would lift the burden currently resting on the researchers' shoulders.
390 For such a task AI needs to be created and trained, and for this training FAIR datasets are needed.
391 This raises the issue that such datasets are simply not available right now in sufficient numbers.

392 As mentioned, data catalogs are among the best tools to implement the FAIR data principles. In
393 order to exploit them to the fullest, it is necessary to standardize the requirements, ensuring that
394 data stored across catalogs is universally findable and accessible. To eventually ensure that data
395 in a data catalog is findable and obtainable without the user actually seeing the user interface of
396 the data catalog. By making the addition of certain (meta)data mandatory, data catalogs can play
397 a vital role in the creation of FAIR metadata sets. With these datasets an AI specialized in tagging
398 data with FAIR metadata can be created and trained. Eventually, with a large enough number of
399 FAIR datasets, this AI would be able to annotate every dataset in a FAIR way, completely lifting
400 this burden off the researchers' shoulders.

401