

# h5RDMtoolbox - A Python Toolbox for FAIR Data Management around HDF5

Matthias Probst <sup>1</sup>Balazs Pritz <sup>1</sup>

1. Institute for Thermal Turbomachinery, Karlsruhe Institute of Technology, Karlsruhe.

**Date Submitted:**

2023-09-26

**Licenses:**This article is licensed under: **Keywords:**

Data management, HDF5, metadata, data lifecycle, Python, database

**Data availability:****Software availability:**Software can be found [here](#)

**Abstract.** Sustainable data management is fundamental to efficient and successful scientific research. The FAIR principles (Findable, Accessible, Interoperable and Reusable) have been proven to be successful guidelines to enable comprehensible analysis, discovery and re-use. Although the topic has recently gained increasing awareness in both academia and industry, the engineering sciences in particular are lagging behind in managing the valuable asset of data. While large collaborations and research facilities have already implemented metadata strategies, smaller research groups and institutes are often missing a common strategy due to heterogeneous and rapidly changing environments as well as missing capacity or expertise. This paper presents an open source package called *h5rdmtoolbox*, written in Python. It is a general-purpose interface to HDF5 files with the aim of helping to quickly implement and maintain FAIR research data management throughout the data lifecycle, using HDF5 as the core file format. One of the key features of the toolbox is the flexible, high-level implementation of metadata standards, adaptable to the changing requirements of projects, collaborations and environments, such as experimental or computational setups. Implementation of interfaces to existing metadata schemas such as *EngMeta* or the *cf-conventions* are possible and part of the comprehensive documentation. Other benefits of the toolbox include a simplified interface to the repository and database solutions for querying metadata stored in HDF5 files.

## 1 Introduction

1 Sustainable data management is fundamental in today's data-driven world for several reasons.  
2 The amount of acquired data storage capacity has long ceased to be the limiting factor, while the  
3 computing power has increased greatly [1]. However, it is the ability to share data rather than  
4 generate it that defines success [2]. Furthermore, interdisciplinary and international collaborations  
5 have become essential in scientific research, and the main means of communication is based on  
6 digital documents [3]. A bottleneck in data exploration and processing, and therefore the general  
7 re-usability, is often the lack of auxiliary data (metadata). As a consequence, much time is spent  
8 on recovering missing information. In some cases, this may require to re-conduct simulations  
9 and experiments. Effective data management practices hence hold the potential of saving time  
10 and money as well as increasing the value of data at the same time.  
11

12 Introducing a new data management concept can be challenging due to conflicting priorities,  
13 expectations, and existing practices, as well as a lack of expertise or clear understanding of  
14 the benefits. Efficient use of standards is crucial for large and interdisciplinary collaborations.  
15 While those groups have developed domains-specific solutions, small research groups and PhD  
16 projects face challenges due to the use of multiple file formats, individual software solutions,  
17 personal preferences for storage and tools, and established structures [4]. Common issues are  
18 the lack of time and resources to develop and implement a comprehensive and sustainable data  
19 management approach [5], which fulfills the requirements of the community and good scientific  
20 practice. Therefore, flexible and manageable solutions are needed to address this issue.

21 Although the implementation of a common management system is beneficial in the long term,  
22 both financially [6] and in terms of efficiency, it disrupts structures and requires time, resources  
23 and cultural change. In academia, high staff turnover is an additional barrier, making it difficult  
24 to establish sustainable solutions. The decay of value develops as projects progress, ultimately  
25 finish and contracts expire. Consequently, the value of data will diminish over time. This issue  
26 is discussed in more detail in [7], [8]. In addition, a value decay can also be observed with  
27 increasing distance from the source of the data. The further away and therefore less involved  
28 a potential data user is, the more information may be missing, either due to restricted access  
29 or limited personal connections. Ensuring that data is preserved and being interpretable at all  
30 times can be achieved by adhering to the so-called FAIR principles, which stand for Findable,  
31 Accessible, Interoperable and Reusable and were first introduced in 2016 by [9]. Since their  
32 publication, the principles have become the cornerstones of many scientific communities and  
33 help to establish a sustainable data management [10]. Structured, highly descriptive information  
34 about data, known as metadata, is an integral part of it. Metadata provides context about its  
35 creation, purpose, use, processing history and the meaning of datasets. Consequently, it enables  
36 data to be discoverable, interoperable and reusable.

37 This work is a contribution to assist small collaborative groups or communities and doctoral  
38 researchers with achieving a FAIR research data lifecycle by using the HDF5 file format. These  
39 groups are often faced with challenges such as heterogeneous file formats, the absence of  
40 standards within their fields, and limited expertise and resources for sustainable data management.  
41 The paper describes the scope and concepts of a Python package named *h5rdmtoolbox* and how it  
42 facilitates the implementation of FAIR principles using the HDF5 file format. Complementing this  
43 manuscript, an extensive online documentation is provided [11], leveraging Jupyter Notebooks  
44 [12]. This documentation offers in-depth insights and additional examples for immediate usage,  
45 serving as a comprehensive resource for users seeking detailed information and practical guidance.

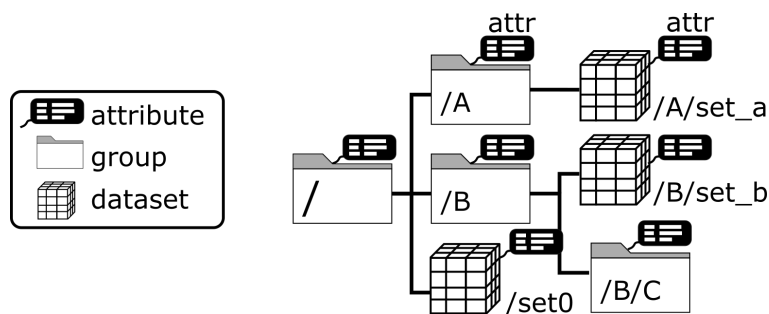
#### 46 **1.1 Outline of the paper**

47 Firstly, the paper outlines the package's scope in comparison to existing and related works. This is  
48 followed by a section stating the concepts and architecture of the toolbox, describing the applied  
49 design principles and methods. Subsequently, the paper discusses concrete implementation  
50 details of all sub-packages and provides illustrative examples, referencing to their relevance  
51 within the research data lifecycle. Limitations of the presented package are stated before the  
52 paper concludes and summarizes the presented work. An outlook is given on future developments  
53 and potential enhancements.

## 54 2 Scope and related work

55 The primary aim of this toolbox is to offer comprehensive support throughout the lifecycle of  
 56 research data (c.f. Figure 2) for small collaborative groups, communities, and doctoral researchers  
 57 engaged in utilizing or contemplating the use of HDF5 files as their central file format. The file  
 58 format is selected for various reasons, which are stated hereafter. A review of other file formats  
 59 is beyond the scope of this work and literature should be referred to, for example [2], [8], [13].

60 HDF5 features efficient storing of large multidimensional datasets together with metadata inde-  
 61 pendent of the storage media, programming environment or operating system. The hierarchical  
 62 structure of group and dataset objects (cf. Figure 1) resembles most engineering data. Attributes  
 63 (key-value pairs) are means to store metadata and can be assigned to each object. The HDF5 file  
 64 format is therefore regarded as self-explanatory. HDF5 finds application in numerous scientific  
 65 domains, such as earth observation [14], high-energy physics [15] or fluid dynamics [16]. An  
 66 in-depth presentation of the file format can be found in [17].



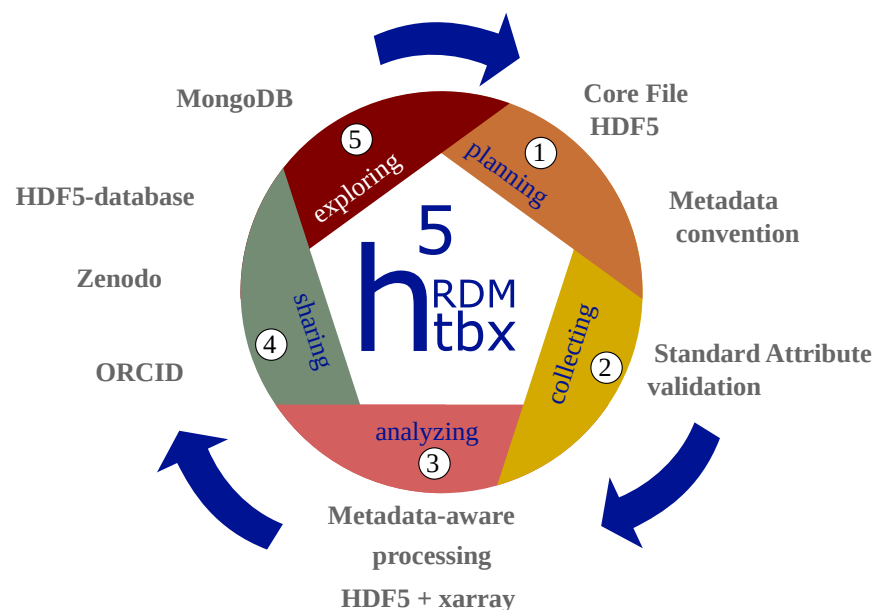
**Figure 1:** Illustration of the hierarchical structure of an HDF file. The internal file structure is organized like a file storage system, where folders are represented by the HDF group objects and files by HDF dataset objects. Both objects can be associated with attributes, which provide the metadata in order to make the objects interpretable.

67 Despite all the advantages of the file format, the organization of data management around HDF5  
 68 is left to the user [18]. This means that the choice of attribute names and values is not regulated by  
 69 any standard. Findability, effective re-usability and automatic analysis, however, are dependent  
 70 on standardization [19].

71 The necessity for designing management solutions around the HDF5 file format is therefore  
 72 evident. While existing solutions, such as [8], [14]–[16], [20], [21] address this need, they  
 73 are often domain-specific, primarily focused on efficiently meeting the demands of specific  
 74 communities rather than providing a generalized framework applicable to diverse problems. For  
 75 example, formats like *Nexus* [15] or *Photon-HDF5* [21] prescribe specific group and dataset  
 76 organizations and metadata usage tailored to their respective data sources, such as neutron and  
 77 X-ray data and molecule spectroscopy experiments, respectively. Other libraries like *Zarr* [20]  
 78 address challenges associated with very large data (terabyte-scale) in the field of bioimaging  
 79 with a particular emphasis on optimized cloud-based operations and the sharing of HDF5-based  
 80 datasets. Finally, questions of efficient database solutions for HDF5 are addressed in [1] and  
 81 [22].

82 Besides the specificity of the solutions, adopting aforementioned solutions to new problems is

83 very difficult due to their complexity and required expertise in the field. When data management  
 84 solutions are needed for a concrete projects, it is crucial to minimize entry barriers. Currently, for  
 85 HDF5, a general approach to manage data in all aspects during its lifecycle including metadata  
 86 concepts, database solutions and practical interfaces are missing. The presented Python package  
 87 *h5rdmtoolbox* seeks to bridge the gap between the advanced communities with domain-specific  
 88 solutions and researchers trying to manage their data without established standards in place.  
 89 Leveraging well-established Python packages, this toolbox offers high-level tools and interfaces  
 90 within one package, that actively contribute to the promotion of FAIR data creation. As a whole,  
 91 the package seeks to be a central resource of tools for scientists allowing them to manage their  
 92 HDF5 data along the full data lifecycle from planning (1) via acquisition (2) and analysis to  
 93 publication in data repositories (4) and sharing in databases (5). Figure 2 illustrates these stages  
 94 and relates keywords to features of the toolbox.



**Figure 2:** Illustration of the lifecycle of research data. Each phase is supported by the *h5rdmtoolbox*. It starts by selecting a file format (HDF5) and a metadata concept (1) and performing quality assurance measures during the selection and processing phase (2). Data is analyzed effectively for scientific output in the next step (3). After publication, the availability of the data should be ensured (4). (Meta)data quality finally is defined by its findability and consequently its re-usable (5) for additional analysis at later time. The respective tools and solutions provided by the toolbox and explained in this word are indicated by keywords around the lifecycle.

### 95 3 Concepts and architecture of the toolbox

96 A key aspect of the toolbox lies in its adaptable implementation of metadata standards and  
 97 interfaces to databases and repositories, allowing it to be relevant across many research fields  
 98 with varying requirements. The challenge is to attain this flexibility without introducing excessive  
 99 complexity, all while ensuring adherence to the FAIR principles. The toolbox achieves this  
 100 through four principles:

101 1. **Relevant programming language:** The choice of programming language significantly

102 impacts the usability and acceptance of this toolbox, as well as data handling in general.  
103 Python is selected for this purpose due to its status as one of the most popular and widely  
104 used language in the scientific community. The high relevance of Python in the field  
105 allows the toolbox to address as many users as possible.

106 2. **One core file format:** The Hierarchical Data Format (HDF5) [23] is selected as the core  
107 and general purpose file format. It is suitable for most scientific and engineering data  
108 sources and allows metadata to be stored with the raw data, making it a self-explanatory  
109 data store. The file format is open-source, well-supported by the HDF5 Group [23] and has  
110 a proven track record in many disciplines. Opting for a single file format as the foundation  
111 for a management toolbox is, therefore, not limiting. Prioritizing user-friendliness and  
112 widespread acceptance, the toolbox implements high-level interfaces to HDF5, extending  
113 the capabilities of the commonly used Python package *h5py*. [24].

114 3. **Flexible Metadata Standardization:** Enabling the storage of metadata alongside raw  
115 data necessitates its standardization (convention) to achieve discoverability. The toolbox  
116 introduces a simple and flexible definition of so-called standard attributes. Users can design  
117 their own convention, which provides feedback about the correctness of the (meta)data  
118 created. This ensures that users maintain the accuracy and completeness of their data and  
119 metadata.

120 4. **Extensibility:** Adaptability extends beyond just metadata standards; it encompasses  
121 various aspects of the toolbox, including interfaces to databases and data repositories.  
122 Abstract classes establish communication rules between HDF5 and users, enabling the  
123 community to add new interfaces on top of the currently implemented ones and to make  
124 them available to others through the toolbox.

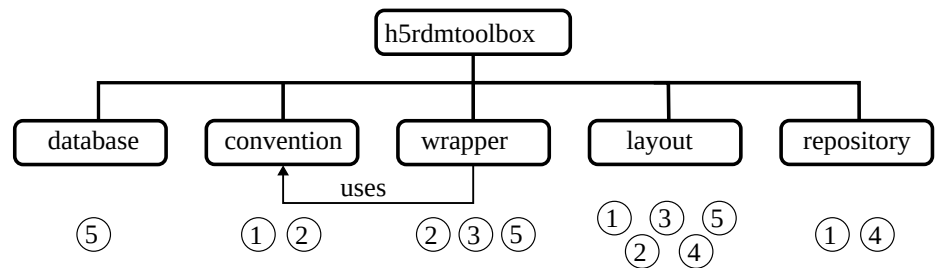
125 In this work, a five-stage representation of the research data lifecycle is adopted, as illustrated  
126 in Figure 2. This framework forms the basis for the toolbox's architectural design, aligning its  
127 functionalities with the key stages of the data lifecycle. Consequently, the toolbox is structured  
128 into five sub-packages and depicted in Figure 3. The numerical assignments in the figure directly  
129 correlate with the roles of these sub-packages in the stages of the data lifecycle (c.f. Figure 2).  
130 This structured approach enhances the toolbox's utility by providing specialized tools for each  
131 phase of the research data lifecycle.

132 Each sub-package corresponds to a key component in the management of HDF5 files throughout  
133 the data lifecycle. These components are designed in a manner that ensures independence from  
134 each other, facilitating individual development and modularity. One exception is made to the  
135 sub-packages *wrapper* and *conventions*.

136 The following sections will highlight the features and implementations of the sub-packages, as  
137 well as their importance within the data lifecycle.

### 138 3.1 layout

139 Research projects start with a scientific question and a data management plan (DMP) [25]. The  
140 DMP outlines how data is handled during and after the project. One important aspect is the  
141 agreement on common exchange formats (in this work HDF5). It has a significant impact on



**Figure 3:** Organization of the sub-packages in the presented package. The core module is called *wrapper*, which adds useful functionality for the user when interacting with HDF5 files. It uses the *convention* module to manage metadata requirements when creating and reading data. The other modules are not dependent on each other and must be imported as required. The numbers indicate their main areas of application within the different stages of the data lifecycle, as shown in [Figure 2](#).

142 the realization of a FAIR data cycle as a whole, especially, when it comes to sharing data [20].  
 143 Besides a common vocabulary, the internal structure (layout) of the file is important. It is the  
 144 basis for reliable processing and automated analysis. The hierarchical structure of HDF5 files  
 145 allows various strategies to organize data and therefore must be regulated by the project data  
 146 manager.

147 The sub-package *layout* implements the class *Layout*, which is a collection of specifications.  
 148 Each specification is a query which will be called during the validation of a file. The motivation  
 149 behind it is, that everything, that is expected to exist in a file, must be findable. The code in  
 150 [Listing 1](#) illustrates the definition of a layout: First, the layout object is created. Then two  
 151 specifications are added by providing a function and query parameters (using pseudocode for  
 152 simplicity). The first parameter can be any function which is able to take the query parameters  
 153 and returns a list of found HDF5 objects. As part of the toolbox and its documentation<sup>1</sup>, a  
 154 database solution provided in the sub-package *database* is used.

```

155 from h5rdmtoolbox import layout
156 lay = layout.Layout()
157 spec1 = lay.add(func=database_function,
158               query=<has dataset with name "x_velocity">)
159 spec2 = lay.add(func=database_function,
160               query=<datasets have the attribute "units">)
161 lay.validate(<hdf5 filename>)
  
```

**Listing 1:** Pseudocode example for a layout definition and validation of an HDF5 file.

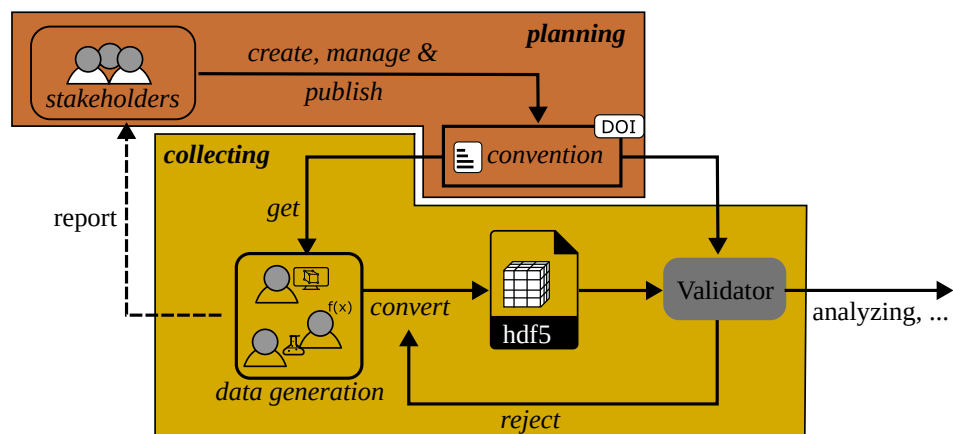
162 The layout concept should be part of every phase in the data lifecycle. After its definition in the  
 163 planning phase (1) it can be useful to validate the integrity of the file at each stage, as the content  
 164 may have been changed. To check if the layout still adheres to the planned definition contributes  
 165 to reliable data and complete files. Avoiding missing information through careful definition of  
 166 the file content in combination with regular checks is the basis of FAIR data.

1. <https://h5rdmtoolbox.rtf.d.io/en/latest/layout/index.html>

### 167 3.2 convention

168 In addition to a robust HDF5 layout, meaningful and comprehensive metadata for HDF5 datasets  
 169 and groups is crucial, ensuring that files are interpretable by both humans and machines. During  
 170 the planning phase (1), a selection of relevant attributes for the investigated problem is important.  
 171 The quality of these attributes significantly influences the findability of data within an HDF5 file,  
 172 as well as the reusability and interoperability aspects of FAIR data in general. The term "quality"  
 173 here refers to whether attributes are linked to existing metadata concepts that can be referenced  
 174 to persistent sources. Examples include controlled vocabularies such as the *cf-conventions* [19],  
 175 metadata schemas like *EngMeta* [26], or ontologies like *metadata4ing* [27]. These sources  
 176 provide standardized and well-defined terms that enhance the clarity and consistency of metadata,  
 177 contributing to improved data discoverability and reuse. Documentation for the toolbox [11]  
 178 includes examples showcasing the possible utilization of these standards within the toolbox. The  
 179 concept of conventions is explained in the following.

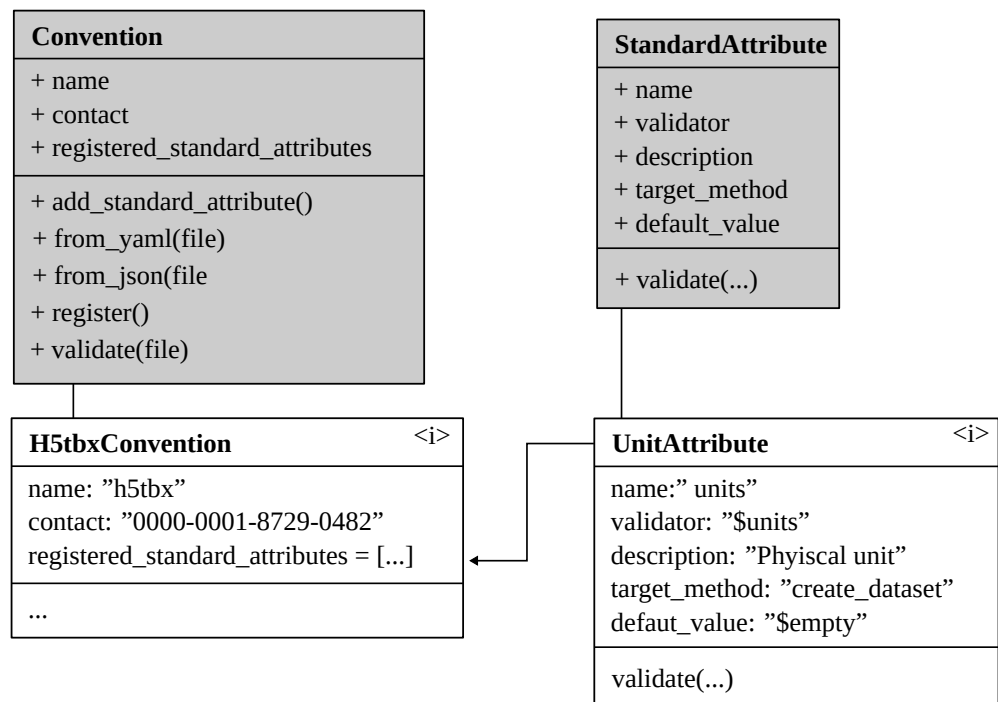
180 The *h5rdmtoolbox* implements the concept of so-called "standard attributes" as part of a "conven-  
 181 tion" to validate relevant metadata (HDF5 attributes) during runtime as the user writes data to the  
 182 file. Figure 4 illustrates a common workflow, which makes use of this concept. The stakeholders  
 183 of a project define and share a set of standardized attributes in the form of a convention. This  
 184 document is shared across all users, which are directly working with HDF5 files. By integrating  
 185 the convention into their workflows through the *h5rdmtoolbox*, they obtain direct feedback  
 186 through a validation mechanism. As a result, the quality in terms of reliable and comprehensive  
 data description through attribute is ensured and basis for the FAIRness of HDF5 file is set.



**Figure 4:** Workflow of collecting and converting the source data. The convention validates the created HDF5 files and serves as a feedback loop to the file creators or the software developers writing the conversion scripts. Only validated files can be further processed or published.

187

188 As shown in the class diagram in Figure 5, a *Convention* objects takes a list of *StandardAttribute*  
 189 objects. The important properties of a *StandardAttribute* are *validator* and *target\_method*. The  
 190 *target\_method* assigns the object to a method of the *h5py* package (other options are `__init__`  
 191 or `create_group`) and the *validator* defines how the attribute is validated during assignment. A  
 192 minimal example of the two instances shown in Figure 5 is written in code in Listing 2. Note,



**Figure 5:** Class diagram of components *Convention* and *StandardAttribute*. The instances "H5tbxConvention" and "UnitAttribute" are used in Listing 2.

193 that the parameter "units" in the function call is not part of the underlying *h5py* package but  
 194 gets dynamically added by enabling the convention. It is also noteworthy that by setting the  
 195 keyword "\$empty" as the default value, the attribute becomes obligatory. For HDF5 datasets,  
 196 this is in fact a reasonable choice, as generally physical data is written to datasets, which require  
 197 a physical unit.

```

198 import h5rdmtoolbox as h5tbx
199 h5tbx.use('h5tbx')
200 with h5tbx.File(myfile.hdf, 'r+') as h5:
201     h5.create_dataset('ds', data=4, units='m/s')
  
```

**Listing 2:** Minimal example of using convention. By enabling the "h5tbx"-convention, the standard attribute "units" becomes obligatory in the method *create\_dataset*. The value of "units" is validated and automatically added to the newly created dataset or an error is raised.

202 As the class diagram suggests, conventions can also be defined in files (JSON or YAML). By  
 203 doing so, the convention can be shared via data repositories or databases with all involved users.  
 204 By enabling the project convention while manipulating the file, users get immediate feedback,  
 205 whether the standardized attributes are valid or not (c.f. Figure 4. This is a difference to the  
 206 concept of layouts, which are static validators. For further information and examples about the  
 207 implementation details, pre-implemented validators as well as the user-defined creation of new  
 208 ones, please refer to the documentation. For a more in-depth understanding of the implementation  
 209 details, pre-implemented validators, and the creation of user-defined validators, please consult



210 the documentation [11], as this information exceeds the scope of this paper. The documentation  
211 provides extensive details, practical examples, and guidance to support users in utilizing and  
212 customizing conventions and validators within the *h5rdmtoolbox*.

### 213 3.3 wrapper

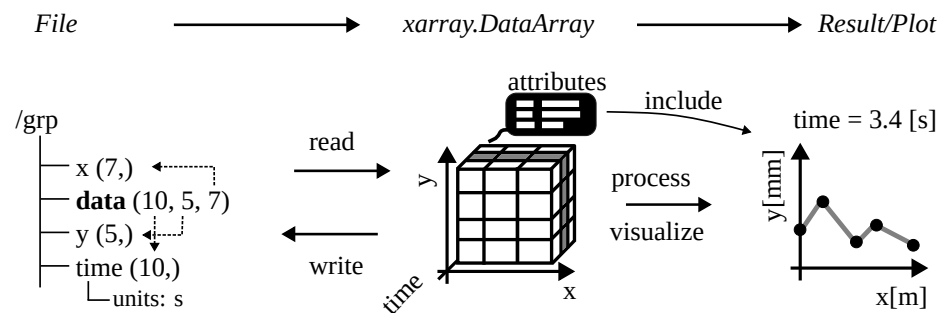
214 The package *wrapper* plays a central role within the toolbox by implementing a thin layer  
215 around the HDF5 Python library *h5py*. Besides user-friendly high-level methods for interactive  
216 representation of the file content in Jupyter Notebooks or helper methods for special datasets,  
217 the *wrapper* package is responsible for

- 218 • integration of the convention concept into the *h5py* framework and
- 219 • metadata-aware exchange of data through *xarray* object [28].

220 The integration of the *xarray* package into the toolbox provides several advantages. As previously  
221 highlighted, one of the reasons for selecting HDF5 is its compatibility with the multidimension-  
222 ality of many scientific and engineering datasets, allowing the storage of attributes alongside  
223 the data. However, using *numpy* arrays as part of the *h5py* package results in the loss of two  
224 important sets of information. Firstly, *numpy* arrays can only represent array data, discarding  
225 attributes associated with the data. Secondly, the axis of a multidimensional array can only  
226 be addressed by their indices (0, 1, etc.), potentially losing references to other datasets in the  
227 HDF5 format (a concept known as *dimension scales* in *h5py* [24]). This limitation hinders the  
228 interpretation of values and their context.

229 The *xarray* package addresses these limitations by wrapping its functionality around *numpy*  
230 arrays [28]. It enables the association of attributes to the values and allows the labeling of the  
231 axes in multidimensional arrays. This structure closely aligns with the HDF5 dataset model. By  
232 returning "metadata-aware" *xarray* objects, the toolbox ensures that provenance information is  
233 added, enhancing the intuitiveness and reliability of data processing. The auxiliary information is  
234 consistently preserved during data utilization for visualization or other post-processing steps, as  
235 depicted in Figure 6. It is noteworthy, that *xarray* has a strong plotting utility, that automatically  
236 extracts information from the data object, incorporating it into the labels and title of the plot.  
237 The synergies between HDF5 and *xarray*, resulting in benefits like concise code and interactive  
238 visualization of metadata, are best illustrated through practical examples. To gain a deeper  
239 understanding and explore enhanced workflows and data operations, it is recommended to  
240 consult the online documentation of the *h5rdmtoolbox*. For the sake of completeness, a short  
241 example is given in the following. The code example in Listing 4 demonstrates the workflow  
242 as illustrated in Figure 6. A subset of the dataset "data" is selected based on the coordinates.  
243 The return value is a *xarray.DataArray* on which the rolling mean is computed. The result is  
244 finally plotted on the screen. With only a few lines of code, the user obtains quick insight into  
245 the dataset while maintaining comprehensibility and traceability.

```
246 import h5rdmtoolbox as h5tbx
247 with h5tbx.File(filename) as h5:
248     # select and read selected data and store in variable:
249     d = h5['data'].sel(x=4.3, y=0.2, method='nearest')
```



**Figure 6:** The h5RDMtoolbox makes use of the *xarray* features. Instead of *numpy* arrays, *xarray.DataArray* objects are returned, which allow carrying the dimension references and attributes and results in comprehensive data processing and visualization.

```

250
251 # process (compute rolling mean over time with window size 3):
252 drm = d.rolling(time=3).mean()
253
254 # visualize the result:
255 drm.plot()

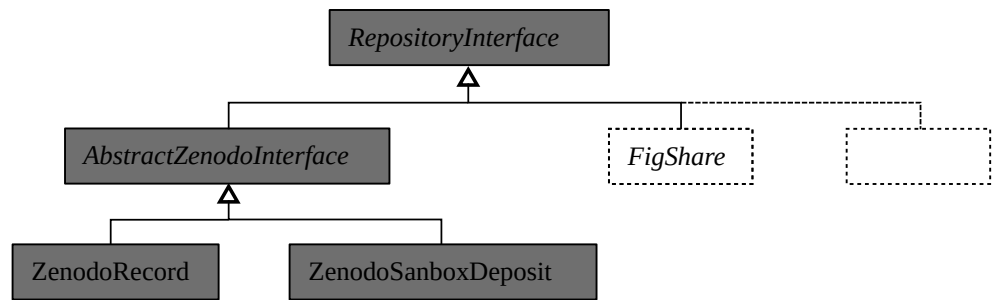
```

**Listing 3:** Example of data extraction using the toolbox. The returned value is an *xarray.DataArray* containing comprehensive metadata from the underlying HDF5 dataset. This facilitates transparent data operations and minimizes potential errors. Additionally, many operations can be reduced to one line of code, which makes scripts concise and traceable.

### 256 3.4 repository

257 How data is shared depends on the scope and restrictions of the project (phase 4 in the lifecycle).  
 258 Most use-cases will, at least for some time, store data locally for internal use or later upload  
 259 to a data repository. The sub-package *repository* implements an abstract interface class from  
 260 which concrete interfaces can be inherited. At the time of writing, only one such repository  
 261 interface is implemented, namely *Zenodo* [29]. As *Zenodo* offers a testing and a production  
 262 interface, two classes are depicted in [Figure 7](#). Although *Zenodo* is one of the most popular  
 263 repositories in the scientific community to publish scientific data which is open-access, more  
 264 platform interfaces are planned to be added. The design of the *repository* sub-package explicitly  
 265 promotes this by using an object-oriented design with an abstract interface class defining the  
 266 user-platform-interaction. [Figure 7](#) depicts the class design. A new repository called *figshare*  
 267 [30] is added in the figure, although not being implemented in the current version. However, it  
 268 illustrates a next possible implementation, which can be contributed by the community.

269 The implemented interface class provides two option of uploaded files. Either the file is uploaded  
 270 as it is, or a metadata file is generated before and then uploaded together with the HDF5 file.  
 271 This has the following reasoning: Large HDF5 files are expensive to download in terms of  
 272 time, especially if it turns out, that the data is not matching the expectations of a user. As data



**Figure 7:** The architecture of the *repository* sub-package allows adding new repository interfaces beyond the existing Zenodo interfaces. The abstract class *RepositoryInterface* sets the rules, such that other popular platforms such as *figshare* could be added by code contributors to the *h5rdmtoolbox*. The gray components are implemented, the dashed lines indicate potential extensions.

273 repositories generally only provide descriptive information about the type of the data publication  
 274 (creator, version, time, keywords, license...) the content of large files can be only be investigated  
 275 after its download. By uploading a metadata file, the internal structure and the attributes can be  
 276 made available in a very small text file (e.g. JSON). It then becomes very feasible to download  
 277 the small metadata file, investigate the content and then decide whether to download to large  
 278 HDF5 file or not. This concept is illustrated in [Figure 8](#).

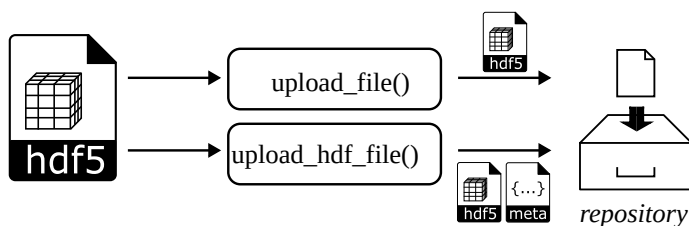
```

279 from h5rdmtoolbox.repository import zenodo
280 repo = zenodo.ZenodoSandboxDeposit(None) # new testing deposit
281
282 # add metadata of data publication, not file content itself:
283 from h5rdmtoolbox.repository.zenodo import metadata
284 repo.metadata = zenodo.metadata.Metadata(...)
285
286 from h5rdmtoolbox.repository.h5metamapper import hdf2json
287 repo.upload_hdf_file('my_file.hdf5', metamapper=hdf2json)
288
289 # download specific filename from the repo:
290 meta_filename = repo.download('my_file.json') # get metadata info first
291 meta_filename = repo.download('my_file.hdf') # if needed
  
```

**Listing 4:** Example of data extraction using the toolbox. The returned value is an *xarray.DataArray* containing comprehensive metadata from the underlying HDF5 dataset. This facilitates transparent data operations and minimizes potential errors. Additionally, many operations can be reduced to one line of code, which makes scripts concise and traceable.

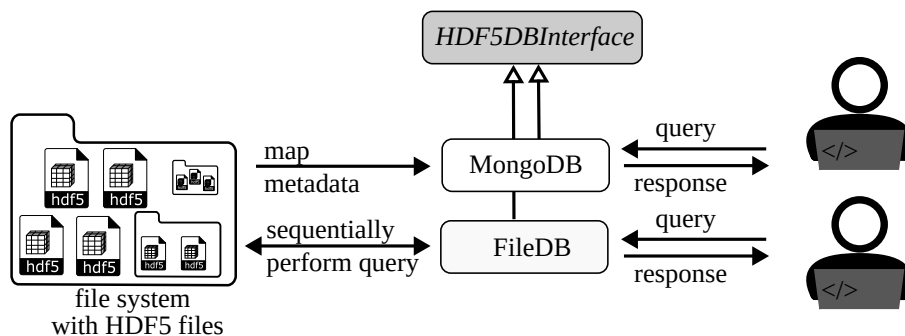
### 292 3.5 database

293 The toolbox implements two ways of sharing and re-using by means of databases:



**Figure 8:** The upload of files to a repository can be performed in two ways: A file can be either uploaded directly or by first deriving a metadata file and then uploading both files. This is especially helpful for large HDF5 files. Interested users may first download the metadata file and inspect the content before downloading the large file.

- 294 1. Using HDF5 as a database inside a file system.
- 295 2. Mapping HDF5 to the NoSQL database MongoDB [31].



**Figure 9:** Workflow of the two provided database solutions: The metadata of the HDF5 files can be mapped to a MongoDB database and then filtered. The other option does not require a database infrastructure but filters the HDF5 files sequentially and returns the data directly. For both solutions, the query functions and return values are the same, because the interfaces are inherited from the abstract database interface class.

- 296 **Figure 9** shows the workflows for both present options. The first approach treats an HDF5 file  
 297 itself as a database and multiple files as multiple databases respectively. The consequences are  
 298 twofold: Firstly, no third party database needs to be used and set up. Secondly, because files  
 299 are opened and searched sequentially, the performance of finding data is not compatible with a  
 300 dedicated system. However, if there are only a few files or the search is within a single file, the  
 301 inefficiency is outweighing. In addition, this concept, as implemented in the toolbox, requires  
 302 no prior operations on the data and only takes a minimum number of lines for the user.
- 303 The second approach extracts the file information and all metadata of each HDF5 object (datasets  
 304 and groups) and writes it into the MongoDB. While this is the most time-consuming part, the  
 305 query itself becomes highly efficient. The used syntax for the queries and the capabilities of the  
 306 present solutions are outlined in detail in the documentation.

### 307 3.6 Documentation

308 The *h5rdmtoolbox* is versioned via a GitHub repository and can be installed using the python  
309 package installer (pip). The current version is v1.0.2 and extensive documentation is automatically  
310 created and published [11]. The documentation website is generated based on Jupyter Notebooks.  
311 On the one hand, this results in practical documentation, showing code and explanations together.  
312 On the other hand, it allows users to reuse the code from the documentation for immediate  
313 application by simply copying the code snippets. As Jupyter Notebooks become more popular  
314 [12], [32], the option to download the full Notebooks will be another efficient option for most  
315 users who are new to the toolbox.

### 316 4 Limitations

317 As outlined before, this package serves as a general toolbox, introducing a management layer  
318 around HDF5 files. Therefore, its strength lies in the metadata organization and user-friendly  
319 interaction with HDF5 files, rather than high-performance data processing. This means, that  
320 during dataset creation and reading, additional processing is needed to validate the metadata  
321 usage. The process of actually writing and reading the file is dependent on the underlying  
322 package, which is *h5py*. For large datasets, however, the overhead is neglectable and the write  
323 or read process is dominating. No significant time differences between *h5py* and the toolbox  
324 are observed. The same accounts for the performance of working with dataset values. They  
325 are provided as *xarray* objects. Again, generating them based on the *numpy* array and other  
326 information from the HDF5 file requires some time. After this, the performance is dependent on  
327 the *xarray* package.

328 The chosen design principles introduce two inherent limitations. Firstly, the package's imple-  
329 mentation in Python inherently limits its compatibility to Python versions. Despite Python's  
330 widespread popularity justifying this choice, it may pose a limitation for users familiar with  
331 other scripting languages like Matlab. While a similar implementation in other languages is  
332 theoretically possible, such an extension is beyond the scope of this work. Secondly, the selection  
333 of HDF5 as the core scientific file format imposes an inherent limitation. Not all scientific or  
334 engineering data may be well-suited for HDF5 files. While HDF5 is versatile, some specialized  
335 data types or structures may find more suitable alternatives outside the HDF5 format.

336 Finally, it is essential to note that the number of interfaces to databases and repositories is  
337 currently limited. As of the current writing, the *database* sub-package includes implementations  
338 for *mongoDB* and a query solution using HDF5 itself. In the *repository* sub-package, only  
339 *Zenodo* is provided. Nevertheless, the toolbox is designed to permit and explicitly encourages  
340 further extensions by the community. This open architecture invites collaborative contributions  
341 to expand the range of interfaces and integrations with databases and repositories based on the  
342 evolving needs and preferences of users.

### 343 5 Conclusion and Outlook

344 The Python package *h5rdmtoolbox* has been introduced, which is designed to support small  
345 collaborative groups, communities, and doctoral researchers who use or consider using HDF5

346 files as their central file format. HDF5 is chosen for its self-descriptive capabilities and versatility  
347 in various scientific domains. However, the management of metadata and internal organization  
348 of datasets and groups, as well as facilitating interoperability with other frameworks, is left  
349 to the users. The toolbox aims to enhance the FAIR principles of data by providing general,  
350 comprehensive tools for managing HDF5 files throughout their lifecycle. While existing solutions  
351 exist to address management needs, they tend to be domain-specific and lack a generalized  
352 framework applicable to diverse problems. Some solutions may only focus on specific aspects of  
353 the data lifecycle, such as databases. In contrast, the presented toolbox adopts a broad approach,  
354 providing tools that enable users to create tailored management solutions for HDF5 files based on  
355 their specific scientific context. Rather than prescribing a singular solution, the toolbox fills the  
356 gap between well-established solutions utilized by large scientific communities and the absence  
357 of standards for individual researchers. By offering a Python package equipped with high-level  
358 tools and interfaces for HDF5 data management, the toolbox improves the FAIRness of HDF5  
359 files for scientists.

360 With user-friendliness and low entry barriers in mind, the toolbox uses popular Python packages  
361 like ‘xarray’ and ‘pydantic’ as dependencies and adopts syntax into newly programmed solutions  
362 (e.g. query within HDF5 files is adopted from mongoDB). However, the toolbox is missing  
363 graphical user interfaces. This would strongly improve the usability and will lower the entry  
364 level, especially for less experienced programmers. Future work should set the focus on the  
365 design of conventions and layouts, as this constitutes the bases of successful data management.

366 Finally, additional testing is required. Although unit tests are implemented, *practical* testing must  
367 identify needs, weaknesses and thus potential for improvements. Application to various problems  
368 and scientific disciplines are planned and feedback from researchers will need to incorporate  
369 into the toolbox. This will extend capabilities, improve the code and allow it to be adapted to the  
370 needs of users. Current concrete use cases investigate fluid problems, such as computational  
371 fluid dynamics simulations and particle image velocity measurements. Lessons learned from  
372 these areas will be incorporated into future publications, while further examples and guidelines  
373 will be continuously added to the online documentation [11].

## 374 6 Acknowledgements

375 The software was developed in-house without any external funding and no conflicts of interested  
376 are declared. The authors would like to thank all users, who have been testing the toolbox so far  
377 and provided helpful feedback. A special thanks belong to Lucas Büttner for the helpful testing  
378 and feedback at the beginning of the project.

## 379 7 Roles and contributions

380 **Matthias Probst:** Conceptualization, Writing, Software Development – original draft

381 **Balazs Pritz:** Project administration, Formal Analysis, Writing - review & editing

382 **References**

- 383 [1] Y. Wang, Y. Su, and G. Agrawal, “Supporting a Light-Weight Data Management Layer  
384 over HDF5,” in *2013 13th IEEE/ACM International Symposium on Cluster, Cloud, and  
385 Grid Computing*, IEEE, 2013, pp. 335–342. DOI: [10.1109/CCGrid.2013.9](https://doi.org/10.1109/CCGrid.2013.9).
- 386 [2] J. Georgieva, V. Gancheva, and M. Goranova, “Scientific Data Formats,” in *Proceedings  
387 of the 9th WSEAS International Conference on Applied Informatics and Communica-  
388 tions*, ser. AIC’09, Moscow, Russia: World Scientific, Engineering Academy, and Society  
389 (WSEAS), 2009, pp. 19–24, ISBN: 9789604741076.
- 390 [3] E. National Academies of Sciences and Medicine, *Open Science by Design: Realizing a  
391 Vision for 21st Century Research*. Washington, DC: The National Academies Press, 2018.  
392 DOI: [10.17226/25116](https://doi.org/10.17226/25116).
- 393 [4] F. De Carlo, D. Gürsoy, F. Marone, *et al.*, “Scientific data exchange: a schema for HDF5-  
394 based storage of raw and analyzed data,” *Journal of synchrotron radiation*, vol. 21, no. 6,  
395 pp. 1224–1230, 2014.
- 396 [5] C. M. Klingner, M. Denker, S. Grün, *et al.*, “Research data management and data sharing  
397 for reproducible research—results of a community survey of the german national research  
398 data infrastructure initiative neuroscience,” *Eneuro*, vol. 10, no. 2, 2023.
- 399 [6] European Commission and Directorate-General for Research and Innovation, *Cost-benefit  
400 analysis for FAIR research data : cost of not having FAIR research data*. Publications  
401 Office, 2019. DOI: [10.2777/02999](https://doi.org/10.2777/02999).
- 402 [7] W. K. Michener, “Meta-information concepts for ecological data management,” *Ecological  
403 Informatics*, vol. 1, no. 1, pp. 3–7, 2006. DOI: [10.1016/j.ecoinf.2005.08.004](https://doi.org/10.1016/j.ecoinf.2005.08.004).
- 404 [8] N. Preuss, G. Staudter, M. Weber, R. Anderl, and P. F. Pelz, “Methods and technologies for  
405 research-and metadata management in collaborative experimental research,” in *Applied  
406 Mechanics and Materials*, Trans Tech Publ, vol. 885, 2018, pp. 170–183.
- 407 [9] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, *et al.*, “The fair guiding principles for  
408 scientific data management and stewardship,” *Scientific Data*, vol. 3, no. 1, p. 160 018,  
409 2016, ISSN: 2052-4463. DOI: [10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).
- 410 [10] A. Jacobsen, R. de Miranda Azevedo, N. Juty, *et al.*, “FAIR Principles: Interpretations  
411 and Implementation Considerations,” *Data Intelligence*, vol. 2, no. 1-2, pp. 10–29, Jan.  
412 2020, ISSN: 2641-435X. DOI: [10.1162/dint\\_r\\_00024](https://doi.org/10.1162/dint_r_00024).
- 413 [11] Probst, Matthias, *HDF5 Research Data Management Toolbox (Documentation)*. [Online].  
414 Available: <https://h5rdmtoolbox.readthedocs.io/en/latest/>, (accessed:  
415 18.12.2023).
- 416 [12] J. M. Perkel, “Why Jupyter is data scientists’ computational notebook of choice,” *Nature*,  
417 vol. 563, no. 7732, pp. 145–147, 2018.
- 418 [13] P. Greenfield, M. Droettboom, and E. Bray, “ASDF: A new data format for astronomy,”  
419 *Astronomy and computing*, vol. 12, pp. 240–251, 2015.
- 420 [14] E. Taaheri and D. Wynne, “An HDF-EOS and data formatting primer for the ECS project,”  
421 Raytheon Company, Tech. Rep., Mar. 2001.

- 422 [15] P. Klosowski, M. Koennecke, J. Tischler, and R. Osborn, “NeXus: A common format for  
423 the exchange of neutron and synchrotron data,” *Physica B: Condensed Matter*, vol. 241,  
424 pp. 151–153, 1997.
- 425 [16] T. Hauser, “Parallel i/o for the cgns system,” in *42nd AIAA Aerospace Sciences Meeting  
426 and Exhibit*, 2004, p. 1088. DOI: [10.2514/6.2004-1088](https://doi.org/10.2514/6.2004-1088).
- 427 [17] S. Koranne and S. Koranne, “Hierarchical data format 5: HDF5,” *Handbook of open  
428 source tools*, pp. 191–200, 2011.
- 429 [18] S. Poirier, A. Buteau, M. Ounsy, *et al.*, “Common Data Model Access: A Unified Layer to  
430 Access Data From Data Analysis Point OF View,” *Icalepcs, Grenoble, October*, 2011.
- 431 [19] J. Gregory, “The CF metadata standard,” *CLIVAR Exchanges*, vol. 8, no. 4, p. 4, 2003.
- 432 [20] J. Moore and S. Kunis, “Zarr: A cloud-optimized storage for interactive access of large  
433 arrays,” in *Proceedings of the Conference on Research Data Infrastructure*, vol. 1, 2023.
- 434 [21] A. Ingargiola, T. Laurence, R. Boutelle, S. Weiss, and X. Michalet, “Photon-HDF5: an open  
435 file format for timestamp-based single-molecule fluorescence experiments,” *Biophysical  
436 journal*, vol. 110, no. 1, pp. 26–33, 2016.
- 437 [22] L. Gosink, J. Shalf, K. Stockinger, K. Wu, and W. Bethel, “HDF5-FastQuery: Accelerat-  
438 ing complex queries on HDF datasets using fast bitmap indices,” in *18th International  
439 Conference on Scientific and Statistical Database Management (SSDBM’06)*, IEEE, 2006,  
440 pp. 149–158.
- 441 [23] The HDF Group, *Hierarchical Data Format, version 5*. [Online]. Available: [https://www  
442 w.hdfgroup.org/HDF5/](https://www.hdfgroup.org/HDF5/), (accessed: 18.12.2023).
- 443 [24] A. Collette, *Python and HDF5*. O’Reilly Media, Inc., 2013, ISBN: 9781449367831.
- 444 [25] A. Salazar, B. Wentzel, S. Schimmler, R. Gläser, S. Hanf, and S. A. Schunk, “How research  
445 data management plans can help in harmonizing open science and approaches in the digital  
446 economy,” *Chemistry—A European Journal*, vol. 29, no. 9, e202202720, 2023.
- 447 [26] B. Schembera and D. Iglezakis, “EngMeta: metadata for computational engineering,”  
448 *International Journal of Metadata, Semantics and Ontologies*, vol. 14, no. 1, pp. 26–38,  
449 2020.
- 450 [27] D. Iglezakis, D. Terzijska, S. Arndt, *et al.*, “Modelling scientific processes with the m4i  
451 ontology,” in *Proceedings of the Conference on Research Data Infrastructure*, vol. 1,  
452 2023. DOI: [/10.52825/cordi.v1i.271](https://doi.org/10.52825/cordi.v1i.271).
- 453 [28] S. Hoyer and J. Hamman, “xarray: ND labeled arrays and datasets in Python,” *Journal of  
454 Open Research Software*, vol. 5, no. 1, 2017.
- 455 [29] M.-A. Sicilia, E. García-Barriocanal, and S. Sánchez-Alonso, “Community curation in  
456 open dataset repositories: Insights from zenodo,” *Procedia Computer Science*, vol. 106,  
457 pp. 54–60, 2017. DOI: [10.1016/j.procs.2017.03.009](https://doi.org/10.1016/j.procs.2017.03.009).
- 458 [30] M. Thelwall and K. Kousha, “Figshare: A universal repository for academic resource  
459 sharing?” *Online Information Review*, vol. 40, no. 3, pp. 333–346, 2016. DOI: [10.1108  
460 /OIR-06-2015-0190](https://doi.org/10.1108/OIR-06-2015-0190).



- 461 [31] K. Chodorow and M. Dirolf, *MongoDB - The Definitive Guide: Powerful and Scalable*  
462 *Data Storage*. O'Reilly, 2010, pp. I–XVII, 1–193, ISBN: 978-1-449-38156-1.
- 463 [32] T. Kluyver, B. Ragan-Kelley, F. Pérez, *et al.*, “Jupyter Notebooks-a publishing format for  
464 reproducible computational workflows.,” *Elpub*, vol. 2016, pp. 87–90, 2016.