

# Metadata Made Easy - Develop and Use Domain Specific Metadata Schemes by following the dmdScheme approach.

## Authors

- **Rainer M Krug**, University of Zürich, Department of Evolutionary Biology and Environmental Studies, Winterthurerstrasse 190, 8057 Zurich, email: [Rainer.Krug@uzh.ch](mailto:Rainer.Krug@uzh.ch), [Rainer@Krugs.de](mailto:Rainer@Krugs.de)
- **Owen L. Petchey**, University of Zürich, Department of Evolutionary Biology and Environmental Studies, Winterthurerstrasse 190, 8057 Zurich, email: [Owen.Petchey@ieu.uzh.ch](mailto:Owen.Petchey@ieu.uzh.ch)

## Corresponding Author

- **Rainer M Krug**, University of Zürich, Department of Evolutionary Biology and Environmental Studies, Winterthurerstrasse 190, 8057 Zurich, email: [Rainer.Krug@uzh.ch](mailto:Rainer.Krug@uzh.ch), [Rainer@Krugs.de](mailto:Rainer@Krugs.de)

# Metadata Made Easy - Develop and Use Domain Specific Metadata Schemes by following the dmdScheme Approach.

Rainer M Krug<sup>a</sup>, Owen L. Petchey<sup>a</sup>

<sup>a</sup>*Department of Evolutionary Biology and Environmental Studies, Winterthurerstrasse 190, 8057 Zurich*

---

## Abstract

1. Metadata plays an essential role in the long term preservation, reuse, and interoperability of data. Nevertheless, creating useful metadata can be sufficiently difficult and weakly-enough incentivised that many datasets may be accompanied by little or no metadata. One key challenge is, therefore, how to make metadata creation easier and more valuable. We present a solution that involves creating domain specific metadata schemes that are as complex as necessary and as simple as possible. These goals are achieved by co-development between a metadata expert and the researchers (i.e. the data creators). The final product is a bespoke metadata scheme into which researchers can enter information (and validate it) via the simplest of interfaces: a web browser application and a spreadsheet.
2. We provide the R package [dmdScheme](#) (Krug & Petchey, 2019a) for creating a template domain specific scheme. We describe how to create a domain specific scheme from this template, including the iterative co-development process, and the simple methods for using the scheme, and simple methods for quality assessment, improvement, and validation.
3. The process of developing a metadata scheme following the outlined approach was successful, resulting in a metadata scheme which is used for the data generated in our research group. The validation quickly identifies

27 forgotten metadata, as well as inconsistent metadata, therefore improv-  
28 ing the quality of the metadata. Multiple output formats are available,  
29 including XML.

- 30 4. Making the provision of metadata easier while also ensuring high quality  
31 must be a priority for data curation initiatives. We show how both  
32 objectives are achieved by close collaboration between metadata experts  
33 and researchers to create domain specific schemes. A near-future priority  
34 is to provide methods to interface domain specific schemes with general  
35 metadata schemes, such as the Ecological Metadata Language, to increase  
36 interoperability.

37 *Keywords:* metadata quality; data curation; archival; long term storage; R  
38 package;

---

## 39 Glossary

- 40 • **analysis** - processing the **analysis-ready data** in order to address the  
41 research question.
- 42 • **analysis-ready data** - data ready for analysis; may be “ready” for a  
43 limited set of analyses. An example would be abundance of each of the  
44 species in a set of communities (e.g. population dynamic data of ecological  
45 communities). (Contrast with **raw data**.)
- 46 • **data deposit package** - a collection of data and **metadata** files deposited  
47 in a long-term repository. This consists at least one data file and the **rich**  
48 **metadata** describing the data file(s) and associated information. May  
49 often contain multiple data files, each with its own **metadata** file.
- 50 • **domain / research domain** - a grouping of e.g. experiments, research,  
51 and / or questions addressed, whose data sets can be described using  
52 **metadata** following one **metadata scheme** which can be regarded as  
53 **rich metadata**. One example is “Experimental Microbial Ecology” for  
54 which the metadata scheme *emeScheme* (Krug & Petchey, 2019b) was  
55 developed. Fields, such as Ecology and Evolutionary Biology, contain  
56 numerous domains.
- 57 • **domain specific metadata scheme** - a **metadata scheme** for a **do-**  
58 **main**.
- 59 • **FAIR data principles** - guiding principles to make data Findable, Ac-  
60 cessible, Interoperable, and Reusable (Wilkinson et al., 2016).
- 61 • **field specific metadata scheme** - a **metadata scheme** general and  
62 broad enough to apply to an entire field. E.g. the Ecological Metadata  
63 Scheme (EML) (Jones et al., 2019).
- 64 • **long-term storage / preservation** - the process of having data stored  
65 / preserved and accessible for the long-term (i.e. greater than 20 years

66 envisaged).

- 67 • **long-term (storage) repositories** - repositories which offer **long term**  
68 **storage**. Examples are ([‘Zenodo - Research. Shared.’ 2020](#)) or ([‘GBIF,’](#)  
69 [2020](#)). The Zenodo repository currently has plans defined for at least 20  
70 years of operation.
- 71 • **metadata** - data about data. Metadata can be as little as the name of a  
72 variable/column in a spreadsheet of data, though such limited metadata  
73 would likely not be considered **rich metadata**, and may not make the  
74 data FAIR.
- 75 • **metadata scheme** - a formalised description of the **metadata** to be  
76 included in e.g. a **data deposit package**, their formats, and which ones  
77 are compulsory or not. A formal scheme assists with the indexing of the  
78 **metadata** that is required for programmatic searching and extracting  
79 **metadata** and data from repositories.
- 80 • **pre-processing** - the preparation of the **raw data** to make it **analysis-**  
81 **ready**. This should be done by a script to make the process reproducible,  
82 and may use different parameters/methods which need to be adjusted  
83 based on the research question and the **raw-data**.
- 84 • **raw data** - data as provided by the measuring device. This could be  
85 images or videos taken from a camera, tables as returned from machines  
86 or hand-written records.
- 87 • **rich metadata** - defined by the Research Data Alliance ([Research Data](#)  
88 [Alliance, 2017](#)) as “data with enough accurate and relevant attributes to  
89 make it easily findable”.

## 90 Introduction

91 To define a kind of gold standard for data handling, [Wilkinson et al. \(2016\)](#)  
 92 developed the so called FAIR data principles. These define principles to make  
 93 the data **F**indable, **A**ccessible, **I**nteroperable and **R**eusable, and help to assess  
 94 data handling workflows in regards to openness.

95 There are multiple reasons why data should be widely reusable ([Heaton,](#)  
 96 [2008](#); [Bishop & Kuula-Luumi, 2017](#); [Pasquetto, Randles, & Borgman, 2017](#)).  
 97 Widely reusable means that anyone making reasonable efforts could reuse the  
 98 data, and that this would be the case even if the data creator(s) are unavailable.  
 99 “Anyone” includes the creator(s) of the data, other members of the creating  
 100 research group, and any other researcher. Use cases include using data from  
 101 previous experiments to plan new ones, re-analysing data using different or new  
 102 pre-processing or analytical approaches to either compare different methodologies  
 103 ([Dufour & Richard, 2019](#)) or to address new scientific aspects (e.g. the use of  
 104 trait databases [Schneider et al. \(2019\)](#)), meta-analyses (e.g. [A. S. Zimmerman,](#)  
 105 [2008](#); [Culina, Crowther, Ramakers, Gienapp, & Visser, 2018](#)), reproduction of  
 106 the studies, and use of data for teaching and training (e.g. [Atenas & Havemann](#)  
 107 [\(2015\)](#); or [Henty \(2015\)](#)).

108 Being able to reuse data includes finding it, understanding why it was  
 109 collected and how it was generated, understanding which datasets are which,  
 110 understanding which variables contain what information, and understanding  
 111 relationships among variables (e.g. [Gregory, Groth, Scharnhorst, & Wyatt \(2020\)](#);  
 112 [Gregory, Groth, Cousijn, Scharnhorst, & Wyatt \(2019\)](#); [A. S. Zimmerman \(2008\)](#);  
 113 [A. Zimmerman \(2007\)](#)). All this information should be stored in metadata;  
 114 thus metadata are essential for reuse ([A. Zimmerman, 2007](#); [Gregory, Groth,](#)  
 115 [Scharnhorst, & Wyatt, 2020](#)). Furthermore, interoperability (the I of FAIR)  
 116 requires standardised metadata schemes.

117 Metadata schemes have been developed which aim at providing a standardised  
118 structure and vocabulary to be used when providing the metadata. Examples of  
119 these schemes are Darwin Core ([Darwin Core task group, 2014](#)) and the Ecological  
120 Metadata Language (short EML) ([Jones et al., 2019](#)) in the field of biology /  
121 ecology, or more broadly Dublin Core ([‘Dublin Core,’ 2020](#)). Interoperability  
122 is essential for research that relies on combining different datasets, and is  
123 particularly important for data-based interdisciplinary research as this very often  
124 combines data from different sources.

125 Given such important reasons for accompanying data with appropriate meta-  
126 data, why do numerous datasets recently published not include useful metadata  
127 ([Roche, Kruuk, Lanfear, & Binning, 2015](#))? For example, a search for “ecology”  
128 and type “dataset” on the Zenodo website in late June 2020 returned 998 data  
129 deposits. The first ten deposits returned contained no deposits with metadata  
130 corresponding to any particular metadata scheme. Two deposit contained meta-  
131 data in a README.md file and a csv file, and one contained a manuscript  
132 called metadata. The remainder contained little or no metadata (other than  
133 the column names in the datasets, which were not explained in more detail). It  
134 may not be far from the truth to say that the majority of thus far deposited  
135 datasets, at least on Zenodo in the field of ecology, have so little metadata or  
136 such poor quality metadata that they have very little possibility to be reused, at  
137 least without considerable effort and with potentially enduring uncertainty.

138 To have the metadata available requires the producer of the data to provide  
139 it. Therefore the answer to the question of why many datasets are deposited  
140 without rich metadata is that the data creators have not prioritised creating  
141 rich metadata. While there is some interest and some level of prioritisation  
142 (e.g. [Campbell, Micheli-Campbell, & Udyawer \(2018\)](#) showed that especially  
143 early career researcher are participating in curating and sharing their data and

144 metadata) a critical question that follows is how to motivate the creation and  
 145 deposition of appropriate metadata. There are multiple possible answers; one  
 146 that we focus upon is that creating metadata is not easy, and creating metadata  
 147 that conforms to a specific scheme is daunting and difficult for researchers. These  
 148 schemes are relatively complex, as they are not specific to a research domain  
 149 (see glossary for definition of “research domain”), but rather for a broad field.  
 150 The advantages of being applicable to a broad field of science (e.g. consistent  
 151 search across a range a wider range of domains, standardised property names  
 152 and vocabulary for metadata provision, interoperability) comes with the cost  
 153 of being complex, somewhat complex and rather difficult to understand, which  
 154 could represent a significant barrier to use by research scientists not working in  
 155 the field of metadata development.

156 Our aim was to make the process of creating metadata not only easy, but  
 157 also useful for the researcher that created the data and, if at all possible, a  
 158 quite pleasurable experience to create! We follow the suggestion of [Poisot,](#)  
 159 [Bruneau, Gonzalez, Gravel, & Peres-Neto \(2019\)](#), that **domain specific meta-**  
 160 **data schemes** (small and purpose-built schemes) can be part of the solution to  
 161 make ecological data easier to find and reuse. The example we use to illustrate a  
 162 domain specific metadata scheme is from the research domain we term “Experi-  
 163 mental Microbial Ecology” ([Worsfold, Warren, & Petchey, 2009](#); e.g. [Pennekamp](#)  
 164 [et al., 2017](#); and [Altermatt et al., 2015](#)) (hereafter EME). We illustrate using  
 165 this domain due to our familiarity with it and because the experimental studies  
 166 involved can be quite complex. Many measurements are often taken using differ-  
 167 ent methods. Multiple treatments are often applied. Numerous taxa are often  
 168 involved. Various steps of data processing are required to obtain analysis-ready  
 169 data (for examples see [Pennekamp et al., 2017](#); and [Garnier, Hulot, & Petchey,](#)  
 170 [2020](#)) from the measured raw data. The methods used can create large amounts



171 of data (several terabytes). Therefore EME is a sufficiently complex domain to  
172 be used as an illustration.

173 To prevent the proliferation of a multitude of domain specific metadata  
174 schemes, risking little or no interoperability among domains, each domain specific  
175 scheme should be as much as possible linked formally to standardised metadata  
176 schemes. In a sense, a domain specific metadata scheme could be regarded as an  
177 easy to use, familiar, intuitive and pleasurable interface to a more general and  
178 standardised metadata scheme.

179 Three other features of domain specific metadata schemes can increase motiva-  
180 tion of researchers to use them: co-development, easy of use, and data/metadata  
181 validation. Co-development by metadata experts and researchers in respective  
182 domains ensures that the scheme can be shaped by providing input to identify  
183 essential properties to be included in the metadata, and to exclude non-essential  
184 metadata. The goal then is to create a domain specific metadata scheme that  
185 fits that domain. Co-development not only results in a better product, but the  
186 resulting “ownership” of these schemes by researchers is likely to increase moti-  
187 vation to use them, to advertise them, to provide input for further development,  
188 and to include them in teaching and training.

189 Metadata entry should not be technically difficult, and presumably the easier  
190 the better. To accomplish these design goals, we made a metadata entry system  
191 that includes only a web browser based application and a spreadsheet. The sim-  
192 plicity of these interfaces should keep the additional workload for the researchers  
193 as small as possible. Moreover, these methods of metadata entry can be common  
194 across domains, meaning that it is not necessary to teach or learn a different  
195 tool for each domain. Previously developed applications for easy metadata entry  
196 include Morpho, a data management tool for earth, environmental, and ecological  
197 scientists (<https://knb.ecoinformatics.org/tools/morpho>); it is open source, but

198 is no longer maintained by the original development team.

199 Data and metadata validation can help researchers increase the quality of their  
 200 data and metadata, for example by checking that variables in datasets contain  
 201 the information they should, and that they correspond to the stated experimental  
 202 treatment and observations. Most large metadata schemes provide mechanisms  
 203 for validating the metadata (e.g. EML in the R package EML ([Boettiger &](#)  
 204 [Jones, 2019](#))). These validations assess mainly the syntactical correctness of  
 205 the metadata, e.g. if all required fields are provided and if numerical values are  
 206 in the allowed range (if ranges are specified). More detailed (contextual and  
 207 contentual) validation can be provided for more specific situations or for smaller  
 208 domains of research, i.e. for domain specific metadata schemes.

209 In this paper, we present as a case study the experience and results of our  
 210 research group in developing the EME domain specific metadata scheme. We first  
 211 used the R package `dmdScheme` ([Krug & Petchey, 2019a](#)) to create a template  
 212 domain specific metadata scheme `emeScheme` ([Krug & Petchey, 2019b](#)) and then  
 213 customised the template scheme to create the EME scheme. We end with a  
 214 discussion on how these domain specific metadata schemes can be integrated  
 215 into larger metadata schemes by using the example of EML ([Jones et al., 2019](#)).

216 The content of this article focuses on presenting the approach by which a  
 217 domain specific metadata scheme can be created using the `dmdScheme` ([Krug](#)  
 218 [& Petchey, 2019a](#)) R package, and its advantages in bringing domain specific  
 219 metadata schemes to more domains and to facilitate the provision of rich and  
 220 quality assured metadata. This article is supported by two Vignettes: one  
 221 describes the `dmdScheme` and is aimed at developers of new domain specific  
 222 schemes and at users interested in a more detailed description of the package.  
 223 The other vignette is aimed at users of the `emeScheme`, and could be modified for  
 224 users of other domain specific schemes. Both are included in the Supplementary

225 information of this article; updated versions are within the respective R packages.

## 226 **The template dmdScheme Package**

227 The R package [dmdScheme](#) ([Krug & Petchey, 2019a](#)) forms the core of de-  
228 veloping and using domain specific metadata schemes following the `dmdScheme`  
229 approach. It is normally hidden for the researcher user of the domain spe-  
230 cific metadata schemes and mainly of concern for the actual developer of new  
231 metadata schemes.

232 The package contains all the base functionality needed to develop a new  
233 domain specific metadata scheme. It includes functionality to create a spreadsheet  
234 for entering the domain specific metadata, functionality to read the metadata  
235 from that spreadsheet, basic validation functions, and export functions to xml  
236 and templates needed to implement the export to EML. It is important to note,  
237 that the `dmdScheme` package itself should not be used to enter actual metadata,  
238 as it does only contain an example metadata scheme.

239 How to develop a new scheme and how to use the package is explained in  
240 detail in the accompanying vignette [Develop and Use the dmdScheme](#) which is  
241 included in the supplemental material of this article.

242 A second part of the `dmdScheme` approach is a repository of domain-specific  
243 schemes ([Krug, 2020](#)). [Here](#) any developed domain-specific schemes can be  
244 deposited. The R package `dmdScheme` contains functionality to load the selected  
245 scheme from this repository and installs the accompanying R package in a  
246 temporary library. This arrangement makes it possible to use the scheme not  
247 only together with the R package `dmdScheme`, but also in other programming  
248 languages, if so desired.

## 249 **Creating a domain specific metadata scheme**

### 250 *Creating the emeScheme*

251 The scheme **emeScheme** (Krug & Petchey, 2019b) was developed based on the  
252 **dmdScheme** (Krug & Petchey, 2019a) and is tailored for data from Experimental  
253 Microbial Ecology. The motivation to develop this metadata scheme was born  
254 out of the realisation that for long-term storage and retrieval following the FAIR  
255 data principles, metadata and data format standards are needed to be able to  
256 find and retrieve the data at any later stage and to be able to reuse it, even  
257 in the own research environment. Therefore it was decided to develop a rich  
258 metadata scheme which would provide enough metadata to be able to find the  
259 data and to re-use it.

260 An open exchange between the researchers and a programmer developing the  
261 scheme was essential in turning the **emeScheme** into a domain specific metadata  
262 scheme which will be used by researchers to create their metadata. Researchers  
263 were involved in the process of developing the **emeScheme** from the beginning.  
264 This included regular meetings to identify properties in the scheme which are  
265 missing, redundant, or not needed. Finally, the researchers were the first testers  
266 of the metadata scheme.

267 The iterative process involved the following steps:

- 268 1. Defining the objectives for the scheme. This included the objective of FAIR  
269 compliance, but also ease of use and validation functionality.
- 270 2. Development of a first version of the scheme. This was done in a spreadsheet  
271 which was then imported in the package and forms the basis of the scheme  
272 definition.
- 273 3. Pilot of entering data from diverse experiments within the domain. The  
274 diversity of experiments is important, as different experiments require  
275 different metadata properties and even structures.

4. Discussion of experiences of the researchers while entering the metadata, highlighting missing, redundant, or not needed properties in the scheme, etc.
5. Incorporate these into the next revision of the scheme and return to step 2.
6. Finalize the scheme definition package and publish it.

Based on initial discussions, the scheme included information about the experiment itself as well as about the different data sets resulting from different measurements and analysis methods as well as treatments during the experiment. This information about the experiment is valuable contextual metadata. To simplify the provision of the metadata and to avoid duplication of the experimental metadata, all metadata would be entered into one spreadsheet file (with multiple sheets), with any required assignment of metadata to individual datasets done automatically in the final stage of the metadata export.

This iterative process resulted in the spreadsheet `emeScheme.xlsx` ( Fig. 1 and Supplemental Material).

This scheme was then bundled together with additional examples and uploaded to the dmdScheme repository as [emeScheme version 0.9.9](#) ([Krug & Petchey, 2019b](#)).

#### *Enhancing the validation*

Even though the package `dmdScheme` already contains a validation function, the validation is generic and mainly structural. The same applies for the export to `xml`, which only exports to a single `xml` file. Additional functionality in the `emeScheme`, i.e. the contextual and contentual validation and the export of the metadata into one `xml` file per data file, is included in an accompanying `emeScheme` R package ([Krug & Petchey, 2019b](#)).

Validation means the checking of the internal consistency of the metadata, compliance with the allowed and suggested values and types of the metadata as

## Metadata made easy - the dmdScheme approach

The figure displays two screenshots of the 'emeScheme' spreadsheet application. The top screenshot shows the 'Experiment' sheet, and the bottom screenshot shows the 'Species' sheet. Both sheets are part of a metadata file named 'DATA emeScheme v1.0.0'.

**Experiment Sheet:**

propertySet	valueProperty	unit	type	suggestedValues	Description	
Experiment	name		character		The name of the experiment.	ASR-expt1
	temperature		character	treatment, in degrees celsius, measurement	Temperature used for all treatments. If different between treatments, use "treatment" and specify in the Treatment sheet.	20
	light		character	treatment, light, dark, cycle, e.g. 16:8 LD	Light used for all treatments. If different between treatments, use "treatment" and specify in the Treatment sheet.	semi-ambient
	humidity		character	treatment, relative humidity in %	Humidity used for all treatments. If different between treatments, use "treatment" and specify in the Treatment sheet.	ambient
	incubator		character	none, bench	What type of incubator is used.	not given here
	container		character		What type of container is used.	Duran type bottle, red lids, 250ml
	microcosmVolume	ml	numeric		Volume of the microcosm container. <b>Not the volume of the culture medium!</b>	100
	mediaType		character			PPM
	mediaConcentration	g/l	numeric			0.55
	cultureConditions		character	axenic, dirty, clean	Conditions of the cultures for all treatments.	dirty
	communityType		character	treatment, single trophic level, multiple trophic level	Characterisation of the microbe community.	initially unknown
	mediaAdditions		character			Wheat seeds added on specific dates, see file wheat_seed_additions.csv
	duration	days	integer		Length of the experiment in days. <b>This should only include the time in which the measurements were taken!</b>	100
	comment		character		Additional features of the Experiment you want to provide	NA

**Species Sheet:**

propertySet	valueProperty	unit	type	suggestedValues	Description	
Species	speciesID		character		Id of the species and strain. Each speciesID has to be unique.	
	name		character		Scientific name of the species or unknown.	
	strain		character			
	source		character		Where the species was obtained from.	
	density	cells / ml	numeric		Initial density used for all treatments. If different between treatments, use "treatment" and specify in the Treatment sheet.	
	functionalGroup		character		Functional group of the species.	
DATA	tt_1					http://www.lgostandards.atcc.org/products/all/30007.aspx
MULTIPLE ROWS	unknown					
	unknown					
	unknown					
	unknown					

Figure 1: Two example sheets (Experiment and Species) in the emeScheme metadata file of the 'emeScheme' spreadsheet. The complete spreadsheet can be found in the supplemental material 'emeScheme.xlsx'.

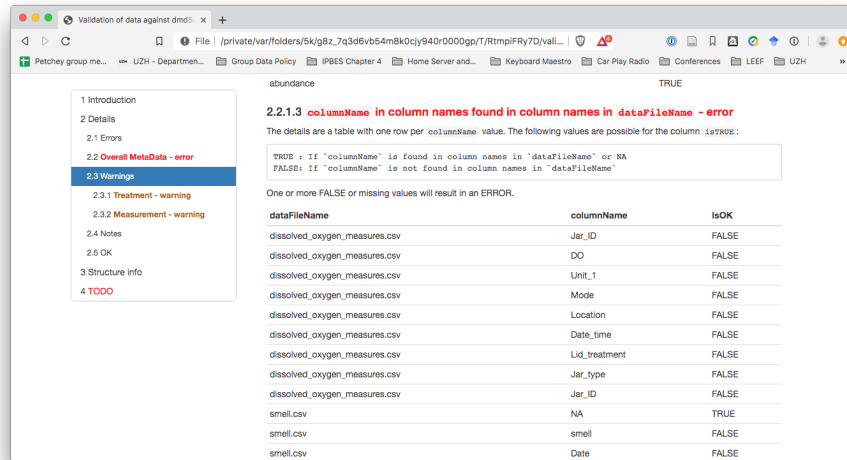


Figure 2: An example of the validation report. The full validation report is in the supplemental material file ‘Validation of data against dmdScheme.html’

well as against the structure of the actual data files. This validation produces an **html** (see Fig. 2), **docx** or **pdf** report, which shows errors, warnings or notes. Errors, warnings, and notes represent different levels of severity of detected faults or inconsistencies in the metadata. For example, if a value is not in the list of **allowed values**, it will result in an **error**, while if it is not in the list of **suggested values**, a **note** will be produced. The validation in the **emeScheme** package had to go beyond the validation included in the **dmdScheme** package. Therefore, it was necessary to write a new validation function to add the new validation rules, i.e. the validation of the **structural metadata** which concerns the data files and its columns.

When the validation has completed without errors, the metadata can be exported to one **xml** file per data file. As in the package **dmdScheme** the export to **xml** creates a single **xml** file, and we needed one **xml** file per data file, a new export function was included in the accompanying R package.

### 317 *Using the emeScheme*

318 The functionality in the **emeScheme**, actually of all **dmdScheme** derived meta-  
 319 data schemes, can be accessed by any of three approaches. As the scheme (and  
 320 the accompanying R package) can be uploaded to the [scheme repository](#) (Krug,  
 321 2020), they are usable from a [universal web app](#) (Krug, 2019) (Fig. 3). Each  
 322 time the web app is started, it re-loads a list of available scheme packages (and  
 323 their accompanying R packages), and these can then be used in the app.

324 Even though this approach is the easiest, it requires the uploading of the  
 325 metadata as well as the data to the server for validation. This might not be  
 326 feasible because of confidentiality / privacy reasons or because of the large size  
 327 of the data files. In this case, the app can also be launched from a local R  
 328 session. The app then runs on the local computer and data never leaves the  
 329 local computer.

330 As a third option, the **emeScheme** and all **dmdScheme** derived packages can  
 331 also be used from the R command line.

332 The different approaches of how this can be done are explained in detail in  
 333 the supplemental material *Develop and Use the dmdScheme*. In addition, the  
 334 document includes detailed information on how new schemes can be developed.  
 335 A more hands on user oriented working example of the **emeScheme** is in the  
 336 supplemental material *emeScheme User Manual*.

### 337 **Integration of Domain Specific Metadata Schemes into the EML Land-** 338 **scape**

339 As mentioned in the Introduction section, interoperability across domains  
 340 requires common cross-domain metadata languages. The **dmdScheme** package  
 341 contains the basic structures to provide an export to EML **xml** format. One of  
 342 the basic requirements of doing so is linking of the domain specific metadata



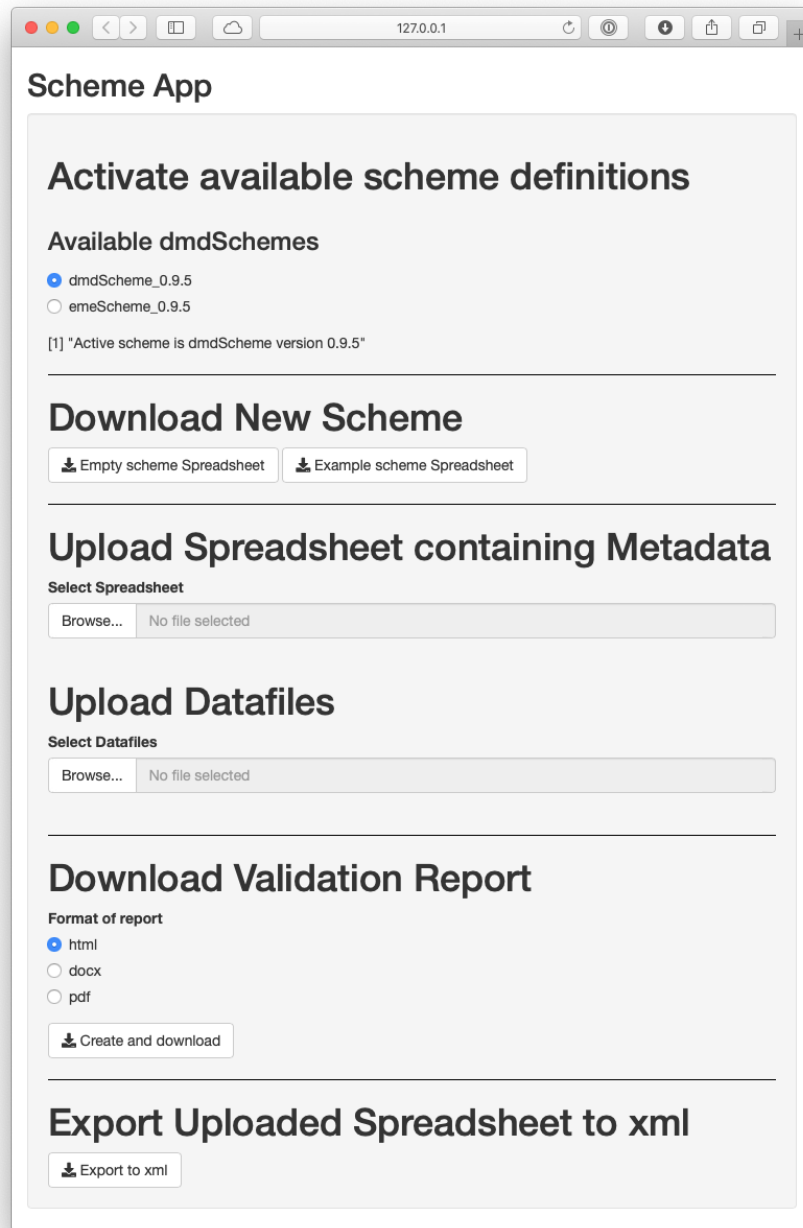


Figure 3: Web app to use the functionality of 'dmdScheme' derived metadata schemes. This app can be run as a [universal web app](#) or also locally.

343 properties to the EML properties. Hence, close inspection of the **emeScheme**  
344 ([Krug & Petchey, 2019b](#)) and some additional constraints (i.e. only one mea-  
345 surement and extraction method per datafile), make it possible to translate the  
346 **emeScheme** metadata into EML. The export into EML is planned for the next  
347 major release of the **emeScheme** package.

348     Export to EML opens a new use case of the **dmdScheme**: if during the  
349 development of a new **dmdScheme** the EML scheme is kept in mind, it will be  
350 possible to use all the functionalities of the package **dmdScheme** as a frontend  
351 for providing EML compliant metadata. In the same way, other large metadata  
352 schemes could be used as the framework for the domain specific metadata  
353 schemes. This would bridge the gap between simple to understand domain  
354 specific metadata schemes on the one side and complex and difficult to understand  
355 but applicable to a large range of different domains metadata schemes.

**Authors Contributions**

RMK and OLP conceived the ideas and designed the methodology; RMK implemented the ideas in R; OLP reviewed the workflow and was the first tester of the package; RMK and OLP led the writing of the manuscript. All authors contributed critically to the drafts and gave final approval for publication.

361 **Data Accessibility**

362 The package does not use any data. The code is available as followed:

- 363 • dmdScheme Package: The package is available on github at [https://github.](https://github.com/Exp-Micro-Ecol-Hub/dmdScheme)  
364 [com/Exp-Micro-Ecol-Hub/dmdScheme](https://github.com/Exp-Micro-Ecol-Hub/dmdScheme) . The version used in this paper  
365 (v1.2) has the doi <https://doi.org/10.5281/zenodo.3894237>
- 366 • emeScheme Package: The package is available on github at [https://github.](https://github.com/Exp-Micro-Ecol-Hub/emeScheme)  
367 [com/Exp-Micro-Ecol-Hub/emeScheme](https://github.com/Exp-Micro-Ecol-Hub/emeScheme) . The version used in this paper  
368 (v1.1.7) has the doi <https://doi.org/10.5281/zenodo.4529180>

369 **Acknowledgements**

370       We have to thank all members of the Predictive ecology Group at the  
371 University of Zurich who provided input in the development of the emeScheme  
372 and functioned as guinea pigs in developing and testing this approach.

373       The authors declare that they have no conflict of interest.

374       Funding was provided by the SNF Project 310030\_188431, and University  
375 Research Priority Programme Global Change and Biodiversity.

## References

- Altermatt, F., Fronhofer, E., Garnier, A., Giometto, A., Hammes, F., Klecka, J.,  
 ... Petchey, O. L. (2015). Big answers from small worlds: A user's guide for  
 protist microcosms as a model system in ecology and evolution. *Methods in  
 Ecology and Evolution*, 6(2), 218–231. doi:[10.1111/2041-210X.12312](https://doi.org/10.1111/2041-210X.12312)
- Atenas, J., & Havemann, L. (2015, November). Open Data as Open Educational  
 Resources: Case Studies of Emerging Practice. doi:[10.6084/m9.figshare.1590031.v1](https://doi.org/10.6084/m9.figshare.1590031.v1)
- Bishop, L., & Kuula-Luumi, A. (2017). Revisiting Qualitative Data Reuse: A  
 Decade On. *SAGE Open*, 7(1), 2158244016685136. doi:[10.1177/2158244016685136](https://doi.org/10.1177/2158244016685136)
- Boettiger, C., & Jones, M. B. (2019). *EML: Read and write ecological metadata  
 language files*.
- Campbell, H. A., Micheli-Campbell, M. A., & Udyawer, V. (2018). Early  
 Career Researchers Embrace Data Sharing. *Trends in Ecology & Evolution*.  
 doi:[10.1016/j.tree.2018.11.010](https://doi.org/10.1016/j.tree.2018.11.010)
- Culina, A., Crowther, T. W., Ramakers, J. J. C., Gienapp, P., & Visser, M.  
 E. (2018). How to do meta-analysis of open datasets. *Nature Ecology &  
 Evolution*, 2(7), 1053–1056. doi:[10.1038/s41559-018-0579-2](https://doi.org/10.1038/s41559-018-0579-2)
- Darwin Core task group, B. I. S. (TDWG). (2014, November). Darwin Core:  
 2014-11-08. Zenodo. doi:[10.5281/ZENODO.12694](https://doi.org/10.5281/ZENODO.12694)
- Dublin Core: Metadata Terms. (2020). <https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>.
- Dufour, I. F., & Richard, M.-C. (2019). Theorizing from secondary qualitative  
 data: A comparison of two data analysis methods. *Cogent Education*, 6(1),  
 1690265. doi:[10.1080/2331186X.2019.1690265](https://doi.org/10.1080/2331186X.2019.1690265)
- Garnier, A., Hulot, F. D., & Petchey, O. L. (2020). Manipulating the strength of  
 organismenvironment feedback increases nonlinearity and apparent hysteresis  
 of ecosystem response to environmental change. *Ecology and Evolution*,

- 10(12), 5527–5543. doi:[10.1002/ece3.6294](https://doi.org/10.1002/ece3.6294)
- GBIF. (2020). <https://www.gbif.org/>.
- Gregory, K., Groth, P., Cousijn, H., Scharnhorst, A., & Wyatt, S. (2019). Searching Data: A Review of Observational Data Retrieval Practices in Selected Disciplines. *Journal of the Association for Information Science and Technology*, 70(5), 419–432. doi:[10.1002/asi.24165](https://doi.org/10.1002/asi.24165)
- Gregory, K., Groth, P., Scharnhorst, A., & Wyatt, S. (2020). Lost or Found? Discovering Data Needed for Research. *Harvard Data Science Review*. doi:[10.1162/99608f92.e38165eb](https://doi.org/10.1162/99608f92.e38165eb)
- Heaton, J. (2008). Secondary Analysis of Qualitative Data: An Overview. *Historical Social Research / Historische Sozialforschung*, 33(3 (125)), 33–45.
- Henty, M. (2015). *Teaching with Research Data* (Report to the {{Australian National Data Service}} ({{ANDS}})) (p. 37).
- Jones, M., O’Brien, M., Mecum, B., Boettiger, C., Schildhauer, M., Maier, M., ... Chong, S. (2019). Ecological metadata language version 2.2.0. doi:[10.5063/f11834t2](https://doi.org/10.5063/f11834t2)
- Krug, R. M. (2019). dmdScheme App. [https://rmkrug.shinyapps.io/dmd\\_app/](https://rmkrug.shinyapps.io/dmd_app/).
- Krug, R. M. (2020, January). dmdSchemeRepository. <https://github.com/Exp-Micro-Ecol-Hub/dmdSchemeRepository>.
- Krug, R. M., & Petchey, O. L. (2019a). dmdScheme: An r package for working with domain specific MetaData schemes (Version v0.9.22). doi:[10.5281/zenodo.3581970](https://doi.org/10.5281/zenodo.3581970)
- Krug, R. M., & Petchey, O. L. (2019b). *emeScheme: A package for working with ecological microbial experimental MetaData*. Zenodo. doi:[10.5281/zenodo.1544945](https://doi.org/10.5281/zenodo.1544945)
- Pasquetto, I., Randles, B., & Borgman, C. (2017). On the Reuse of Scientific Data. *Data Science Journal*, 16(0), 8. doi:[10.5334/dsj-2017-008](https://doi.org/10.5334/dsj-2017-008)
- Pennekamp, F., Griffiths, J. I., Fronhofer, E. A., Garnier, A., Seymour, M., Altermatt, F., & Petchey, O. L. (2017). Dynamic species classification of

- microorganisms across time, abiotic and biotic environmentsA sliding window  
approach. *PLOS ONE*, 12(5), e0176682. doi:[10.1371/journal.pone.0176682](https://doi.org/10.1371/journal.pone.0176682)
- Poisot, T., Bruneau, A., Gonzalez, A., Gravel, D., & Peres-Neto, P. (2019).  
Ecological Data Should Not Be So Hard to Find and Reuse. *Trends in  
Ecology & Evolution*, 0(0). doi:[10.1016/j.tree.2019.04.005](https://doi.org/10.1016/j.tree.2019.04.005)
- Research Data Alliance. (2017). Rich Metadata - DFT. [https://smw-rda.esc.rzg.mpg.de/index.php/Rich\\_Metadata](https://smw-rda.esc.rzg.mpg.de/index.php/Rich_Metadata)
- Roche, D. G., Kruuk, L. E. B., Lanfear, R., & Binning, S. A. (2015). Public  
Data Archiving in Ecology and Evolution: How Well Are We Doing? *PLOS  
Biology*, 13(11), e1002295. doi:[10.1371/journal.pbio.1002295](https://doi.org/10.1371/journal.pbio.1002295)
- Schneider, F. D., Fichtmueller, D., Gossner, M. M., Güntsch, A., Jochum,  
M., König-Ries, B., ... Simons, N. K. (2019). Towards an ecological  
trait-data standard. *Methods in Ecology and Evolution*, 10(12), 2006–2019.  
doi:[10.1111/2041-210X.13288](https://doi.org/10.1111/2041-210X.13288)
- Wilkinson, M. D., Dumontier, M., Aalbersberg, Ij. J., Appleton, G., Axton,  
M., Baak, A., ... Mons, B. (2016). The FAIR Guiding Principles for  
scientific data management and stewardship. *Scientific Data*, 3, 160018.  
doi:[10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18)
- Worsfold, N. T., Warren, P. H., & Petchey, O. L. (2009). Context-dependent  
effects of predator removal from experimental microcosm communities. *Oikos*,  
118(9), 1319–1326. doi:[10.1111/j.1600-0706.2009.17500.x](https://doi.org/10.1111/j.1600-0706.2009.17500.x)
- Zenodo - Research. Shared. (2020). <https://zenodo.org/>.
- Zimmerman, A. (2007). Not by metadata alone: The use of diverse forms of  
knowledge to locate data for reuse. *International Journal on Digital Libraries*,  
7(1-2), 5–16. doi:[10.1007/s00799-007-0015-8](https://doi.org/10.1007/s00799-007-0015-8)
- Zimmerman, A. S. (2008). New Knowledge from Old Data: The Role of  
Standards in the Sharing and Reuse of Ecological Data. *Science, Technology,  
& Human Values*, 33(5), 631–652. doi:[10.1177/0162243907306704](https://doi.org/10.1177/0162243907306704)