



BNL-223193-2022-JAAM

Enabling FAIR data in Earth and environmental science with community-centric (meta)data reporting formats

R. Crystal-Ornelas, A. Rogers

To be published in "Scientific Data"

July 2022

Environmental and Climate Sciences Department
Brookhaven National Laboratory

U.S. Department of Energy

USDOE Office of Science (SC), Biological and Environmental Research (BER) (SC-23)

Notice: This manuscript has been authored by employees of Brookhaven Science Associates, LLC under Contract No. DE-SC0012704 with the U.S. Department of Energy. The publisher by accepting the manuscript for publication acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes.

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or any third party's use or the results of such use of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof or its contractors or subcontractors. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

1 Enabling FAIR data in Earth and environmental science with community-centric (meta)data
2 reporting formats

3

4 **Robert Crystal-Ornelas^{1*}, Charuleka Varadharajan¹, Dylan O’Ryan¹, Kathleen Beilsmith²,**
5 **Ben Bond-Lamberty³, Kristin Boye⁴, Madison Burrus¹, Shreyas Cholia⁵, Danielle S.**
6 **Christianson⁵, Michael Crow⁶, Joan Damerow¹, Kim S. Ely⁷, Amy E. Goldman⁸, Susan**
7 **Heinz⁶, Valerie C. Hendrix⁵, Zarine Kakalia¹, Kayla Mathes⁹, Fianna O’Brien⁵, Stephanie**
8 **Pennington³, Emily Robles¹, Alistair Rogers⁷, Maegen Simmonds^{1,10}, Terri Velliquette⁶,**
9 **Pamela Weisenhorn², Jessica Nicole Welch⁶, Karen Whitenack¹, Deborah A. Agarwal⁵**

10 ¹Earth and Environmental Sciences Area, Lawrence Berkeley National Laboratory, Berkeley, CA
11 94720, USA

12 ²Argonne National Laboratory, Lemont, IL 60439, USA

13 ³Pacific Northwest National Laboratory, Joint Global Change Research Institute at the University
14 of Maryland–College Park, College Park, MD 20740, USA

15 ⁴Geochemistry and Biogeochemistry Group, SLAC National Accelerator Laboratory, 2575 Sand
16 Hill Road, Menlo Park, CA 94025, USA

17 ⁵Scientific Data Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720 USA

18 ⁶Environmental Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN 37830, USA

19 ⁷Environmental and Climate Sciences Department, Brookhaven National Laboratory, Upton, NY
20 11973, USA

21 ⁸Pacific Northwest National Laboratory, Richland, WA 99354, USA

22 ⁹Integrated Life Sciences, Virginia Commonwealth University, Richmond, VA 23284, USA

23 ¹⁰Now at: Pivot Bio, 2910 Seventh Street, Berkeley, CA 94710, USA

24

25

26 Corresponding author: Charuleka Varadharajan (cvaradharajan@lbl.gov)

27 Abstract

28 Research can be more transparent and collaborative by using Findable, Accessible,
29 Interoperable, and Reusable (FAIR) principles to publish Earth and environmental science data.
30 Reporting formats—instructions, templates, and tools for consistently formatting data within a
31 discipline—can help make data more accessible and reusable. However, the immense diversity
32 of data types across Earth science disciplines makes development and adoption challenging.
33 Here, we describe 11 community reporting formats for a diverse set of Earth science (meta)data
34 including cross-domain metadata (dataset metadata, location metadata, sample metadata), file-
35 formatting guidelines (file-level metadata, CSV files, terrestrial model data archiving), and
36 domain-specific reporting formats for some biological, geochemical, and hydrological data
37 (amplicon abundance tables, leaf-level gas exchange, soil respiration, water and sediment
38 chemistry, sensor-based hydrologic measurements). More broadly, we provide guidelines that
39 communities can use to create new (meta)data formats that integrate with their scientific
40 workflows. Such reporting formats have the potential to accelerate scientific discovery and
41 predictions by making it easier for data contributors to provide (meta)data that are more
42 interoperable and reusable.

43 **Keywords:** data standards, metadata, FAIR data, open science, data repository

44 1. Introduction

45 Making Earth and environmental science data Findable, Accessible, Interoperable, and
46 Reusable (FAIR)^{1,2} contributes to research that is more transparent and reproducible³. Search
47 engines and data repositories^{2,4,5} have enabled advances in data preservation, findability, and
48 accessibility. However, data interoperability and reuse remain major challenges in part due to
49 the diversity of Earth science data, and because researchers may lack time and funding for data
50 management or awareness of tools and resources to make data more reusable^{5,6}. This results
51 in barriers to scientific research and knowledge generation; for example, synthesis of data
52 across different sources can be extremely time-consuming when data and metadata are not
53 standardized in a common, well-defined format.

54 Standards for data and metadata, hereafter referred to as (meta)data standards, have
55 been proposed as important elements to make Earth and environmental science data easier to
56 find, understand and reuse⁷⁻¹⁰. Formal (meta)data standards are typically accredited by large
57 governing bodies and emphasize making data broadly reusable¹¹. For example, the
58 International Organization for Standardization (ISO) 8601 standard provides guidelines for
59 formatting date and timestamps and has been adopted in a wide range of research and
60 business sectors¹². The Open Geospatial Consortium's Sensor Observation Service standard¹³
61 outlines standardized ways of pulling sensor data from web interfaces. Such accredited
62 standards are extraordinarily useful, but are available only for a few environmental data types
63 and can take over a decade to build governing processes and consensus.

64 In contrast, reporting formats are community efforts aimed at harmonizing diverse
65 environmental data types without the oversight of the governing protocols or working groups
66 that maintain vocabularies and extensive documentation. There are reporting formats for
67 different research domains and data types including water quality¹⁴ and meteorological data¹⁵.
68 Reporting formats are typically more focused within scientific domains—for example, marine
69 observations¹⁶ or geoscience¹⁷. Reporting formats can enable efficient collection and
70 harmonization of information needed to understand and reuse specific types of data within a
71 research community. For example, the use of FLUXNET's half-hourly flux and meteorological
72 reporting format¹⁸ enables both access and reuse of consistently formatted carbon, water, and
73 energy flux data from thousands of sampling locations across the world. However, reporting
74 formats do not exist for most environmental data types, and even if they do, complexity and lack
75 of resources can limit their adoption⁹.

76 There are many scientific benefits when research communities adopt reporting formats,
77 ranging from organizing data collection in the field or lab to more efficient data reuse in

78 synthesis and modeling efforts. Reporting formats can facilitate data sharing within a group,
79 provide guidelines for consistent data collection, enable streamlined scientific workflows, and
80 enable long-term preservation of knowledge that may not be typically stored or reported with the
81 data^{19,20}. Moreover, research disciplines are beginning to operationalize and implement
82 practices^{21,22} to achieve the original FAIR guiding principles^{21,22}. Reporting formats developed by
83 the research communities for which they are intended are seen as a critical step toward
84 achieving greater data interoperability and reuse²².

85 A variety of multidisciplinary data are generated in research sponsored by the U.S.
86 Department of Energy (DOE) and stored in the Environmental Systems Science Data
87 Infrastructure for a Virtual Ecosystem (ESS-DIVE) data repository^{4,23}. Integration and analysis of
88 diverse data types such as hydrological, geological, ecological, biological, and climatological
89 data is an essential element of complex environmental systems science (ESS) research.
90 However, such interdisciplinary data integration presents unique challenges, such as
91 inconsistent use of terms, formats, and metadata across disciplines²⁴. In this manuscript, we
92 describe and harmonize 11 diverse and complementary (meta)data reporting formats that our
93 interdisciplinary team developed for commonly used data types in ESS research to enable their
94 archival following FAIR principles in general purpose repositories such as ESS-DIVE. These
95 include guidelines to format and describe general research elements (e.g., general file
96 metadata, tabular data, physical samples, model data), as well as guidelines developed for
97 more specific data types relevant to interdisciplinary research (e.g., biogeochemical samples,
98 soil respiration, leaf-level gas exchange). As part of this process, we adopted or used
99 components of existing reporting formats or standards to the greatest extent possible, and also
100 developed new reporting formats for some data types. These can be used individually or
101 collectively in scientific workflows, and many of the formats are widely applicable for
102 environmental research. Moreover, the process we used for developing the formats—including
103 our approach to obtain community consensus, mirror documentation across several web
104 platforms, and track community feedback—can be used by other research communities to
105 develop reporting formats for their own purposes.

106 2. Results

107 Our community-centric approach to developing reporting formats had four key outcomes
108 that are broadly important to making scientific data more reusable. First, the teams reviewed a
109 total of 112 pre-existing data standards and other data resources (data repositories, data
110 systems, datasets, projects) to create (meta)data crosswalks (Supplementary Files 1-20). Such

111 crosswalks provide a tabular map of existing resources related to each data type, allowing the
112 teams to identify gaps in existing standards, and determine which variables, terms, and
113 metadata were essential to harmonize and incorporate into their reporting formats. At the onset
114 of the review process, ESS-DIVE recommended adopting existing standards to the extent
115 possible. However, we found that for all 11 data types, none entirely met ESS research
116 community needs, and this necessitated development of all 11 reporting formats.

117 Second, we created 11 reporting formats (Table 1) that encompass a range of complex
118 and diverse ESS (meta)data fields that can be used when researchers upload data to ESS-
119 DIVE. Six of the reporting formats created by our community of scientists are cross-domain
120 reporting formats (Figure 1a), which apply broadly to data across different scientific disciplines.
121 These reporting formats were developed to help researchers more consistently format their
122 (meta)data for interdisciplinary science applications and include basic dataset metadata for
123 citation and findability²⁵, file-level metadata²⁶, guidelines for formatting Comma Separated Value
124 (CSV) files²⁷, sample metadata²⁸, terrestrial model data archiving guidelines²⁹, and research
125 locations metadata³⁰. The remaining five reporting formats apply to different domain data types
126 (Figure 1b) and include microbial amplicon abundance tables³¹, leaf-level gas exchange³², soil
127 respiration³³, sample-based water and soil chemistry measurements³⁴, and water level and
128 sonde-based hydrologic measurements³⁵. All reporting formats have a minimal set of required
129 metadata fields necessary for programmatic data parsing and optional fields that provide
130 detailed spatial/temporal context about the sample useful to downstream scientific analyses.
131 Throughout development, we aimed to strike a balance between pragmatism for the scientists
132 reporting data and machine-actionability that is emblematic of FAIR data. A comparison
133 between FAIR guiding principles and our reporting formats (Table 3) highlights how a
134 community-centric effort like ours can move data archiving towards achieving many FAIR data
135 principles (though see discussion for limitations).

136 Together, these 11 reporting formats are part of a flexible, modular, and integrated
137 framework (Figure 1) that can accommodate new reporting formats in the future, and enable
138 their findability and accessibility individually or collectively. As part of the framework
139 development, all teams created templates with harmonized terms and formats to be internally
140 consistent as much as possible. For example, dates are always reported in YYYY-MM-DD
141 format. Whenever reporting formats include spatial data, the variables are harmonized as
142 “latitude” and “longitude” and reported in decimal degrees with common bounds (-90 to 90 and -
143 180 to 180, respectively). All formats that require CSV files adopted as many recommendations
144 from the CSV reporting format as possible. Data collected using the water and soil chemistry,

145 and amplicon reporting formats have an option to report a persistent identifier for associated
146 samples [International General Sample Number (IGSN)], to enable effective tracking across
147 online data systems, as outlined in the Sample ID reporting format.

148 The third outcome is related to how we shared and archived all reporting formats in three
149 ways, each with a distinct use. First, all reporting formats are published as datasets in the ESS-
150 DIVE repository, which enables direct, public download and citation upon use. Second, each
151 reporting format is hosted on the version control platform GitHub, which enables ongoing edits
152 and versioning while also allowing users to provide feedback³⁶. Third, the most up-to-date
153 reporting format content from GitHub is rendered as a project website through the service
154 GitBook³⁷. We mirrored the reporting format instructions and templates across several web
155 platforms to ensure the documentation is available in a variety of digital formats to serve the
156 needs of various user groups and stakeholders. GitHub is likely a more familiar platform and
157 user interface for software engineers and informatics specialists, for example, while GitBook
158 websites may be preferred by Earth science researchers.

159 Lastly, we formulated guidelines (Box 1) for research communities that want to replicate
160 our model of community-centric (meta)data reporting format development. We encourage (1)
161 reviewing existing standards, (2) developing a crosswalk of terms across relevant standards or
162 ontologies of interest, (3) iteratively developing templates and documentation with feedback
163 from prospective users, (4) assembling a minimum set of (meta)data required for reuse, and (5)
164 hosting finalized documentation on platforms that can be publicly accessed and updated easily.

165 3. Discussion

166 Many scientific journals and funders require data deposition in long-term repositories.
167 However, in many cases, data are submitted to repositories in bespoke formats with little
168 (meta)data standardization⁵. Community-led (meta)data reporting formats like the set described
169 in this paper can enable archived data to be more reusable and interoperable^{21,22}. Our scientist-
170 centric approach to creating the formats helped to determine workflows that are most useful and
171 practical for researchers to adopt. Here we discuss important aspects that need to be
172 considered in development and use of such reporting formats.

173 Reporting formats can help researchers organize and synthesize their own (meta)data
174 for their research purposes. It can be challenging for small teams, or even individuals to keep
175 track of data collected over multi-year field campaigns or laboratory experiments^{19,20}. Early
176 adoption of a consistent way of compiling data can help individuals or research teams avoid ad

177 hoc data collection practices and also help researchers efficiently integrate their data,
178 particularly when multiple analyses or teams are involved.

179 Moreover, community reporting formats can lead to greater data accessibility and reuse.
180 For example, researchers in the Ameriflux network³⁸ organize flux data in the Flux Processing
181 (FP-in) reporting format¹⁸. When participants in the network agree to provide their flux data in
182 this format³⁹, benefits include: 1) access to data services such as automated QA/QC of datasets
183 and value-added ONEFlux data processing⁴⁰, 2) Digital Object Identifier assignment which helps
184 to track dataset citation and reuse, and 3) potential to increase findability of their data. Similarly,
185 when contributors upload datasets on ESS-DIVE, they are offered automated metadata quality
186 assessments, and published data are assigned DOIs and made searchable across the
187 DataONE network. In another example, the Watershed Function Scientific Focus Area project⁴¹
188 adopted ESS-DIVE's water and soil chemistry reporting format as an initial step towards
189 establishing a field data workflow in a community observatory where diverse hydrological,
190 geochemical, geophysical, ecological, and remote sensing datasets are collected⁴². The use of
191 the reporting format will make it possible for researchers to synthesize data on chemical
192 concentrations both within and across field locations.

193 Application of the reporting formats also allows for the use of tools and services that
194 enhance data curation, findability and reuse. As an example, some of the fields in ESS-DIVE's
195 dataset metadata reporting format²⁵ allow programmatic metadata quality validation, which
196 checks for field presence, format, and length. Because these metadata can be mapped to a
197 variety of machine-readable metadata formats including JSON-LD and the U.S. Department of
198 Energy's Office of Scientific and Technical Information (OSTI) reporting formats⁴³. This further
199 enabled transforming and disseminating ESS-DIVE datasets across other platforms such as
200 Google Dataset Search, DataONE, OSTI and DataCite.

201 The development of these reporting formats was driven by the scientific need for
202 practical tools for data management, while improving the potential for data reuse achieving
203 many of the FAIR guiding principles (Table 3). We made several pragmatic choices to ensure
204 that the reporting formats would have a low barrier to adoption by time-limited researchers. This
205 included investigating whether using pre-existing reporting formats "off the shelf" would meet
206 project and researcher's scientific needs and workflows. Although it is desirable to use existing
207 formats whenever possible, we found that there were many circumstances when they do not
208 directly apply to a scientific community's research (meta)data needs. For example, although the
209 Water Quality Exchange format¹⁴ is used within the United States to report water quality
210 monitoring data by local, state, and federal agencies, the format was not entirely suitable for

211 ESS-DIVE's purposes. Some of the concerns raised by the community included: 1) the structure
212 of the data and metadata templates that are used for regulatory reporting were considered to be
213 cumbersome and inefficient for scientific use (e.g., containing redundant elements of sampling
214 and analytical methodology along with the data) and 2) the required vocabularies (as specified
215 in the template dictionary) were found to be difficult to use because they included several terms
216 that were unnecessary, while missing terms for specific analytes of interest to the community.

217 To address these concerns, we developed the ESS-DIVE reporting format for sample-
218 based water and soil chemistry³⁴ that is more suitable for files typically generated in scientific
219 laboratories. It borrows elements from the WQX standard, but provides flexibility in format and
220 terminology, while capturing sufficient metadata and vocabularies to enable data exploration
221 and reuse including the ability to use scripts to compare and combine different datasets³⁴. In this
222 way, the water and soil chemistry reporting format achieves some component of FAIR guiding
223 principle "I2" that suggests using ontologies, while still being responsive to a research
224 community that desired flexibility in research terminology (Table 3).

225 Similarly, when creating the sample ID metadata reporting format, we decided to extend
226 the existing IGSN sample identifier template and guidelines in ESS-DIVE's Sample ID reporting
227 format to meet researchers' need to link interdisciplinary environmental and biological samples,
228 and to minimize effort in providing information for sample collections⁴⁴. In this case,
229 incorporating IGSNs ensures that researchers using this format achieve FAIR principle "F3" and
230 have globally unique identifiers for their data products, which facilitates tracking associated
231 sample data across multiple online data systems. In an effort to be pragmatic, we decided to
232 lower the threshold for adoption of the sample ID reporting format (and nearly all others; Table
233 3) by compromising on elements that would achieve FAIR principle "I3" related to machine
234 readable knowledge representation. All reporting formats encourage users to define variables in
235 a data dictionary. Though this may not be fully machine readable according to the FAIR
236 principles²¹, this method of defining variables is a key step toward reusable and machine
237 actionable data. The feedback gathered when creating our Sample ID reporting format was then
238 provided to the broader IGSN community to help improve the IGSN metadata template for
239 interdisciplinary science^{45,46}.

240 Through the process, we learned that many (meta)data standards are not accessible to
241 a typical researcher and require a significant learning curve to become fluent in the informatics
242 terminology used by established data standards. For example, the Open Geospatial
243 Consortium's data standard for environmental sensors¹³ is a detailed schema described over
244 100 pages, which is challenging for a typical scientific researcher to understand and implement.

245 Hence, we had to make several pragmatic choices to ensure that the reporting formats would be
246 amenable to adoption by time-limited researchers. Once choice involved replacing terms in
247 existing standards with words that were more intuitive to scientists. For example, whilst there
248 was no reporting format for leaf-level gas exchange data, a crosswalk of the instrument output
249 from a relatively small number of instrument manufacturers quickly identified a common
250 terminology that already had broad acceptance and use by the scientific community
251 (Supplementary File 7). By using crosswalks (Supplementary Files 1-10) our teams were able to
252 map ESS-DIVE's reporting formats to many existing (meta)data standards and other data
253 resources, and, in the future, will allow building tools that enable interoperability with different
254 systems. We also simplified the reporting format templates and instructions to the greatest
255 extent possible by specifying a few required fields and several more optional fields to provide
256 additional details.

257 Our model and guidelines of supporting and empowering the scientific community to
258 develop (meta)data reporting formats that meet their needs can enable other communities to
259 undertake these internal data standardization efforts that make their data even more useful
260 beyond the purpose for which they were collected (Box 1). We acknowledge that other research
261 infrastructures have made important strides toward data standardization within research
262 communities though they can still take dozens of years to manifest¹⁷. We found value in
263 including a broad range of stakeholders in the process, and included field personnel who make
264 the measurements, instrument manufacturers, and scientists who use the data in models or
265 synthesis activities⁴⁷.

266 There are incentives that can help promote widespread adoption of these or other
267 formats to justify the time investment required for individual researchers or teams into scientific
268 workflows. First, involving data collectors and reusers at the core of the development process
269 makes the resulting formats more pragmatic and scientifically useful. Importantly, the domain
270 scientists involved in the reporting format development became community ambassadors and
271 helped engage their use by fellow researchers through conference presentations and peer-
272 reviewed papers^{44,47-49}. Second, we expanded our user community by sharing information about
273 the reporting formats through a series of webinars, documentation, tutorials, and personalized
274 community outreach. These incentives have had some success, as evidenced by the datasets
275 submitted to ESS-DIVE using one or more of the reporting formats within a few months after
276 they were finalized (Table 2).

277 We identify some future work that can potentially lower the barrier to adopting reporting
278 formats, provide added benefits to those who use the formats, and make (meta)data FAIRer¹⁰.

279 Currently, ESS-DIVE applies a set of manual checks to datasets uploaded to ESS-DIVE that
280 follow the reporting format. However, development of automated formatting checkers⁵⁰ would
281 help users instantly validate their datasets against reporting format guidelines. Other types of
282 software can also be built around the reporting formats. For example, software could be
283 developed to automatically convert sensor or instrument-derived data into the units requested
284 by a reporting format. As a starting point for this work, the file-level metadata reporting format
285 already includes an open-source script⁵¹ that enables reading and parsing data files submitted
286 in that format. The leaf-level gas exchange reporting format includes a detailed translation table
287 matching the reporting format data variables with standard outputs from 10 commonly used,
288 commercially produced instruments. This could provide the foundation for development of
289 conversion software to automatically format data with the recommended variable names and
290 units. ESS-DIVE is also planning a data integration and fusion component of the repository that
291 will facilitate synthesizing and analyzing datasets that adhere to any of the 11 ESS-DIVE
292 reporting formats. Enabling advanced queries within the files will require development of
293 software and data parsers so that a great number of reporting formats achieve FAIR principle
294 “F4” which calls for data to be fully searchable.

295 With more data being generated than ever, reusable data can have substantial societal,
296 economic, and scientific impacts. But for Earth and environmental science data, which are
297 complex and heterogeneous, achieving reusability will require concentrated effort at (meta)data
298 standardization within research communities. Our work to develop 11 community (meta)data
299 reporting formats is a critical step to making Earth and environmental science data more
300 reusable because we emphasize human readability that is compatible with machine readability.
301 We hope that our model of empowering research communities to self-organize and create their
302 own (meta)data reporting formats will enable other communities to undertake these internal data
303 standardization efforts that make their data even more useful beyond the purpose for which they
304 were collected.

305 4. Methods

306 Earth and environmental science data are complex, multi-scale, and span diverse research
307 domains such as geology, hydrology, climate, ecology, and biology. At ESS-DIVE, we initiated a
308 community-centric model that engaged domain scientists to develop formats for common Earth
309 science data types. The objective was to create formatting guidelines and templates that would
310 gather the minimum but sufficient metadata or data necessary for data interpretation and reuse.

311

312 **4.1 Reviewing existing standards and feedback on drafts**

313 Each team conducted a review of existing standards (Table 1), involving both literature
314 searches and exploring resources from informatics groups (e.g., Research Data Alliance and
315 Earth Science Information Partners) or agencies working with similar data, to identify whether
316 any existing data standards or conventions could be used directly or to inform their reporting
317 format. Based on this review, each team created tabular ‘crosswalks’ (Supplementary Files 1-
318 10) to map related terminology from relevant standards. This process helped identify gaps in
319 existing standards, and determine important elements that had to be present, and variations in
320 terminology used across different standards that required harmonization. For example, some
321 existing standards report date and time under the column name ‘datetime’ while another reports
322 the same information, as ‘ValueDateTime’ (see example of a terminology crosswalk³⁵). Here, we
323 provide a brief narrative of methods for each reporting format with details on existing data
324 standards and other data resources reviewed during reporting format development. For further
325 details on the technical aspects of each reporting format, please refer to ESS-DIVE’s community
326 space on GitHub³⁶ or view the datasets for each reporting format submitted to ESS-DIVE (Table
327 1).

328

329 **4.2 Obtaining community consensus**

330 Each team created instructions and (meta)data templates for their reporting formats. The teams
331 piloted the formats within their research groups and communities to ensure the templates were
332 practical and useful for scientists who collect and reuse data (Figure 2). In total, 247 individuals
333 representing 128 institutions provided input at various stages of the reporting format
334 development process. As the reporting format instructions and templates reached a final stage,
335 they published the “ready-to-use” reporting formats in three locations each with distinct benefits
336 for the end-users: GitHub³⁷, GitBook, and the ESS-DIVE data repository to enable findability
337 and long-term preservation.

338

339 **4.3 Cross-domain reporting formats**

340 **4.3.1 Dataset metadata**

341 The goal for creating the dataset metadata reporting format was to ensure that any dataset
342 submitted to ESS-DIVE would have complete and descriptive metadata to enable its findability
343 and citation upon use. The ESS-DIVE team reviewed machine and human-readable metadata
344 standards including the Ecological Metadata Language⁵² as well as JSON for Linking Data⁵³.

345 The ESS-DIVE metadata reporting format follows existing metadata standards as much as
346 possible (e.g., 'title' in Ecological Metadata Language is also 'title' for ESS-DIVE's metadata).

347

348 **4.3.2 File-level metadata**

349 The file-level metadata reporting format was developed for users to provide details about the
350 individual files contained within a dataset. The review of existing standards²⁶ included file-level
351 metadata used across 6 organizations (e.g., USGS, NEON).

352

353 **4.3.3 CSV file formatting guidelines**

354 The CSV reporting format was developed to provide guidelines for more consistently formatting
355 tabular data²⁷. The intention was to make this a domain agnostic set of guidelines so that
356 anyone who works with tabular data can use the format in their research to make tabular data
357 more interoperable and machine-readable. The team reviewed existing standards and
358 guidelines (Table 1) including recommendations from the Environmental Data Initiative (e.g., do
359 not mix data types in a column) and the ORNL DAAC (e.g., indicating missing numeric values
360 with -9999).

361

362 **4.3.4 Sample IDs and metadata**

363 The ESS-DIVE Sample ID reporting format²⁸ aligns as much as possible with extensive work on
364 IGSN⁵⁴ with the goal of standardizing sample collection metadata and more efficiently tracking
365 physical samples sent to different collaborators, labs, data systems, etc. This work also
366 reviewed 12 different standards and data resources to provide recommendations for improving
367 interoperability of biological^{8,55} and environmental samples¹⁴.

368

369 **4.3.5 Terrestrial model data archiving guidelines**

370 The model data archiving reporting format²⁹ was informed by input from the DOE's land
371 modeling community and other guidelines from the American Geophysical Union and National
372 Science Foundation Earthcube communities. In developing the guidelines⁴⁹, the goal was to
373 help modelers make decisions about which components of their terrestrial models should be
374 archived in a long-term data repository. The guidelines were developed with input on which
375 model data were most useful to archive, how long they remained useful, and what scientific
376 purpose they would serve.

377

378 **4.3.6 Location metadata**

379 The goal of developing the location metadata reporting format was to provide generalized
380 guidelines for describing locations used in research. The review of existing standards included
381 metadata templates from specific projects at some of the DOE's National Labs to understand
382 the different field sampling strategies of large interdisciplinary projects. The review also included
383 known standards and guidelines for recording locations such as Climate and Forecast
384 Conventions⁵⁶, the Federal Geographic Data Committee's Content Standard for Digital
385 Geospatial Metadata⁵⁷ and the Open Geospatial Consortium⁵⁸.

386

387 **4.4 Reporting formats for domain-specific data types**

388 In addition to the set of 6 cross-domain reporting formats described above, we also developed 5
389 formats that are tailored to specific data types commonly used in the terrestrial and subsurface
390 ecosystem research community. ESS-DIVE's goal was to engage Earth and environmental
391 scientists to determine practical reporting formats that data contributors are willing to use while
392 at the same time ensuring a high potential for data reuse.

393

394 **4.4.1 Amplicon abundance table metadata**

395 The reporting format for amplicon abundance table metadata was developed to facilitate
396 consistent reporting of microbiome sample data with the format of these tables following ESS-
397 DIVE's CSV file guidelines. Required data (e.g., representative sequences) were chosen to
398 support comparisons of abundance tables across studies. The reporting format distinguishes
399 between sequencing metadata and bioinformatic processing metadata for amplicon abundance
400 tables. As much as possible, the team aligned recommendations for sequencing metadata with
401 the existing standards developed by the Genomic Standards Consortium for minimum
402 information about a marker gene sequence and minimum information about any (x) sequence⁵⁵
403 (Supplementary File 6). In the absence of an existing standard for bioinformatic processing
404 metadata, the reporting format contains a minimal set of fields to capture the data processing
405 steps most relevant for comparing and combining amplicon counts across studies (Table 1).
406 The final set of sequencing and bioinformatic metadata fields selected were informed by a
407 community of scientists involved with either the development of microbiome data pipelines or
408 conducting microbiome studies in both field and lab settings.

409

410 **4.4.2 Leaf-level gas exchange**

411 The team working on this reporting format³² reviewed existing conventions used in plant trait
412 databases, large data collections developed for synthesis papers, and the variable descriptions

413 that are part of standard instrument outputs in order to determine the most suitable variable
414 names to use to report leaf-level gas exchange data. Templates for formatting metadata about
415 the methods and sample materials used in an experiment, as well as details on the
416 instrumentation involved in collecting data were developed through an iterative process of input
417 and review open to all interested stakeholders. The reporting format is designed to be flexible
418 and modular, provides guidelines on the archive of raw and processed data, and seeks to
419 capture experimental metadata needed to interpret and reuse these data types⁴⁷.

420

421 **4.4.3 Soil respiration**

422 To create the soil respiration reporting format, this team reviewed and integrated
423 recommendations from 9 existing guidelines and standards (Table 1)³³. The review captured an
424 array of how different standards format their general metadata and data (e.g., formatting date
425 and timestamps) and also accounted for a range of soil-atmosphere gas exchange data types
426 (e.g., GHGs or radiocarbon)⁴⁸.

427

428 **4.4.4 Sample-based water and soil chemistry measurements**

429 The goal in creating a reporting format for water-soil-sediment data was to harmonize chemical
430 concentration data that span several measurement types. The review included 15 standards
431 (Table 1) for related environmental chemistry measurements including metadata elements from
432 the EPA's WQX¹⁴ as well as EarthChem⁵⁹. Based on input from the potential ESS user
433 community that included both data collectors, managers, and modelers, we developed a
434 reporting format based on community input³⁴.

435

436 **4.4.5 Water level and sonde-based hydrologic monitoring**

437 This reporting format harmonizes variables common to sonde-based hydrologic monitoring
438 research including water level, temperature, and pH data. The existing standards and/or data
439 sources included in the crosswalk for the hydrologic monitoring reporting format (Table 1) were
440 chosen for inclusion given their common use in the scientific community. They aligned generally
441 on the types of hydrologic metadata to record (e.g., information about dates and times as well
442 as information about data collection sites) but had different terminology across each of the
443 resources³⁵. The development of the reporting format included a review of additional data
444 sources and standards beyond those listed in the crosswalk (Table 1).

445

446 Acknowledgements

447 Robert Crystal-Ornelas, Charuleka Varadharajan, Dylan O’Ryan, Madison Burrus, Shreyas
448 Cholia, Joan Damerow, Valerie C. Hendrix, Zarine Kakalia, Fianna O’Brien, Emily Robles,
449 Maegen Simmonds, Karen Whitenack, and Deborah A. Agarwal were funded through the ESS-
450 DIVE repository by the U.S. DOE’s Office of Science Biological and Environmental Research
451 under contract number DE-AC02-05CH11231. Kim S. Ely and Alistair Rogers were supported
452 through the US Department of Energy contract number DE-SC0012704 to Brookhaven National
453 Laboratory. We acknowledge the work of Diana Swantek in producing the Figure 2 illustration.
454 Reporting format development was supported by ESS-DIVE’s Community Funds, through the
455 Office of Biological and Environmental Research in the Department of Energy, Office of
456 Science.

457 Competing interests

458 The authors of this manuscript have no competing interests to declare.

459 Data availability

460 Each data reporting format and all supporting documentation are hosted on our GitHub
461 Community Space³⁶ and archived in the ESS-DIVE data repository^{25–35}. The supplementary
462 information for this manuscript is also archived in ESS-DIVE⁶⁰.

463 Code availability

464 We have made code available which enables file-level metadata extraction⁵¹ for files that
465 adhere to the reporting format.

466 Author contributions (according to CRediT taxonomy)
467 Conceptualization: DAA, CV
468 Data curation: RCO, CV, DO, KB, BBL, KB, MC, JD, KSE, AEG, SH, KM, SP, AR, MS, TV, PW,
469 JNW, DAA
470 Funding Acquisition: DAA, CV
471 Methodology: DAA, CV, RCO, JED, KB, BBL, KB, MC, JD, KSE, AEG, SH, KM, SP, AR, MS,
472 TV, PW, JNW
473 Project Administration: DAA, KW
474 Resources: DAA
475 Software: MC
476 Supervision: DAA, CV
477 Visualization: RCO, CV
478 Writing - original draft: RCO, CV, JED
479 Writing - review and editing: RCO, CV, DO, KB, BBL, KB, MB, SC, DSC, MC, JD, KSE, AEG,
480 SH, VCH, ZK, KM, FO, SP, ER, AR, MS, TV, PW, JWN, KW, DAA
481

482 **References**

483

- 484 1. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
485 stewardship. *Sci Data* **3**, 160018 (2016).
- 486 2. Stall, S. *et al.* Advancing FAIR data in Earth, space, and environmental science. *Eos, Earth
487 and Space Science News* **99**, (2018).
- 488 3. Toelch, U. & Ostwald, D. Digital open science—Teaching digital tools for reproducible and
489 transparent research. *PLoS Biol.* **16**, e2006022 (2018).
- 490 4. Varadharajan, C. *et al.* Launching an Accessible Archive of Environmental Data. *EOS* **100**,
491 (2019).
- 492 5. Tenopir, C. *et al.* Data sharing, management, use, and reuse: Practices and perceptions of
493 scientists worldwide. *PLoS One* **15**, e0229003 (2020).
- 494 6. Perrier, L., Blondal, E. & MacDonald, H. The views, perspectives, and experiences of
495 academic researchers with data sharing and reuse: A meta-synthesis. *PLoS One* **15**,
496 e0229182 (2020).
- 497 7. Sansone, S. A. *et al.* FAIRsharing as a community approach to standards, repositories and
498 policies. *Nature Biotechnology* **37**, 358–367 (2019).
- 499 8. Wieczorek, J. *et al.* Darwin Core: an evolving community-developed biodiversity data
500 standard. *PLoS One* **7**, e29715 (2012).
- 501 9. Poisot, T., Bruneau, A., Gonzalez, A., Gravel, D. & Peres-Neto, P. Ecological Data Should
502 Not Be So Hard to Find and Reuse. *Trends Ecol. Evol.* **34**, (2019).
- 503 10. Michener, W. K. Ecological data sharing. *Ecol. Inform.* **29**, 33–44 (2015).
- 504 11. Lin, D. *et al.* The TRUST Principles for digital repositories. *Sci Data* **7**, 144 (2020).
- 505 12. ISO. Date and Time Format (ISO Standard Number 8601-1:2019). (2019).
- 506 13. Bröring, A., Stasch, C. & Echterhoff, J. OGC Sensor Observation Service Interface
507 Standard, Version 2.0. (2012).

- 508 14. Read, E. K. *et al.* Water quality data for national-scale aquatic research: The Water Quality
509 Portal. *Water Resour. Res.* **53**, 1735–1745 (2017).
- 510 15. AmeriFlux. *BADM: Biological, Ancillary, Disturbance, and Metadata*
511 <https://ameriflux.lbl.gov/data/badm/> (2020).
- 512 16. Dañobeitia, J. J. *et al.* Toward a Comprehensive and Integrated Strategy of the European
513 Marine Research Infrastructures for Ocean Observations. *Frontiers in Marine Science* **7**,
514 (2020).
- 515 17. Cocco, M. *et al.* The EPOS Research Infrastructure: a federated approach to integrate solid
516 Earth science data and services. *Ann. Geophys.* **65**, DM208–DM208 (2022).
- 517 18. Flux Processing (FP-In). *Half-Hourly / Hourly Data Upload Format*
518 <https://ameriflux.lbl.gov/half-hourly-hourly-data-upload-format/> (2017).
- 519 19. Goodman, A. *et al.* Ten Simple Rules for the Care and Feeding of Scientific Data. *PLoS*
520 *Comput. Biol.* **10**, e1003542 (2014).
- 521 20. Lowndes, J. S. S. *et al.* Our path to better science in less time using open data science
522 tools. *Nat Ecol Evol* **1**, 0160 (2017).
- 523 21. Jacobsen, A. *et al.* FAIR principles: Interpretations and implementation considerations.
524 *Data Intellegence* **2**, 10–29 (2020).
- 525 22. Bailo, D. *et al.* Perspectives on the Implementation of FAIR Principles in Solid Earth
526 Research Infrastructures. *Front Earth Sci. Chin.* **8**, (2020).
- 527 23. Environmental System Science Data Infrastructure for a Virtual Ecosystem. *ESS-DIVE*
528 <https://data.ess-dive.lbl.gov/data>.
- 529 24. Hills, D. J. *et al.* Earth and space science informatics perspectives on integrated,
530 coordinated, open, networked (ICON) science. *Earth Space Sci.* **9**, (2022).
- 531 25. Agarwal, D. *et al.* ESS-DIVE Reporting Format for Dataset Package Metadata. *ESS-DIVE*
532 <https://www.doi.org/10.15485/1866026> (2022).
- 533 26. Velliquette, T. *et al.* ESS-DIVE Reporting Format for File-level Metadata. *ESS-DIVE*

- 534 <https://www.doi.org/10.15485/1734840> (2021).
- 535 27. Velliquette, T. *et al.* ESS-DIVE Reporting Format for Comma-separated Values (CSV) File
536 Structure. *ESS-DIVE* <https://www.doi.org/10.15485/1734841> (2021).
- 537 28. Damerow, J. *et al.* ESS-DIVE global sample numbers and metadata reporting format for
538 Environmental Systems Science (IGSN-ESS). *ESS-DIVE*
539 <https://www.doi.org/10.15485/1660470> (2020).
- 540 29. Simmonds, M. B. *et al.* ESS-DIVE guidelines for archiving terrestrial model data. *ESS-DIVE*
541 <https://www.doi.org/10.15485/1813868> (2021).
- 542 30. Crystal-Ornelas, R. *et al.* ESS-DIVE Reporting Format for Location Metadata. *ESS-DIVE*
543 <https://www.doi.org/10.15485/1865730> (2022).
- 544 31. Weisenhorn, P. & Beilsmith, K. ESS-DIVE Reporting Format for Amplicon Abundance
545 Tables. *ESS-DIVE* <https://www.doi.org/10.15485/1865729> (2022).
- 546 32. Ely, K. S., Rogers, A. & Crystal-Ornelas, R. ESS-DIVE reporting format for leaf-level gas
547 exchange data and metadata. *ESS-DIVE* <https://www.doi.org/10.15485/1659484> (2020).
- 548 33. Bond-Lamberty, B., Christianson, D. S., Crystal-Ornelas, R., Mathes, K. & Pennington, S.
549 C. ESS-DIVE reporting format for field measurements of soil respiration. *ESS-DIVE*
550 <https://www.doi.org/10.15485/1798520> (2021).
- 551 34. Boye, K. *et al.* ESS-DIVE Reporting Format for Sample-based Water and Soil Chemistry
552 Measurements. *ESS-DIVE* <https://www.doi.org/10.15485/1865731> (2022).
- 553 35. Goldman, A. E., Ren, H., Torgeson, J. & Zhou, H. ESS-DIVE Reporting Format for
554 Hydrologic Monitoring Data and Metadata. *ESS-DIVE*
555 <https://www.doi.org/10.15485/1822940> (2021).
- 556 36. ESS-DIVE Community Space. *ESS-DIVE Community Space* [https://github.com/ess-dive-](https://github.com/ess-dive-community)
557 [community](https://github.com/ess-dive-community) (2021).
- 558 37. Crystal-Ornelas, R. *et al.* A guide to using GitHub for developing and versioning data
559 standards and reporting formats. *Earth Space Sci.* **8**, (2021).

- 560 38. Novick, K. A. *et al.* The AmeriFlux network: A coalition of the willing. *Agric. For. Meteorol.*
561 **249**, 444–456 (2018).
- 562 39. AmeriFlux Data Policy. *AmeriFlux Data Policy* <https://ameriflux.lbl.gov/data/data-policy/>
563 (2021).
- 564 40. Onboarding and Orientation for new site teams. [https://ameriflux.lbl.gov/sites/onboarding-](https://ameriflux.lbl.gov/sites/onboarding-and-orientation-for-new-site-teams/)
565 [and-orientation-for-new-site-teams/](https://ameriflux.lbl.gov/sites/onboarding-and-orientation-for-new-site-teams/) (2017).
- 566 41. Hubbard, S. S. *et al.* The East River, Colorado, watershed: A mountainous community
567 testbed for improving predictive understanding of multiscale hydrological–biogeochemical
568 dynamics. *Vadose Zone J.* **17**, 1–25 (2018).
- 569 42. Kakalia, Z. *et al.* The Colorado East River community observatory data collection. *Hydrol.*
570 *Process.* **35**, (2021).
- 571 43. OSTI. *Instructions for announcement of U.S. Department of Energy (DOE) publicly*
572 *available scientific research datasets* <https://www.osti.gov/elink/F2416instruct.jsp> (2017).
- 573 44. Damerow, J. *et al.* Sample identifiers and metadata to support data management and reuse
574 in multidisciplinary ecosystem sciences. *Data Science Journal* **20**, 1–19 (2021).
- 575 45. Klump, J. *et al.* Towards globally unique identification of physical samples: Governance and
576 technical implementation of the IGSN global sample number. *Data Sci. J.* **20**, (2021).
- 577 46. Richard, S. M. *et al.* Internet of samples. *Proc. Assoc. Inf. Sci. Technol.* **58**, 813–815
578 (2021).
- 579 47. Ely, K. S. *et al.* A reporting format for leaf-level gas exchange data and metadata. *Ecol.*
580 *Inform.* (2021).
- 581 48. Bond-Lamberty, B., Christianson, D. S., Crystal-Ornelas, R., Mathes, K. & Pennington, S.
582 C. A reporting format for field measurements of soil respiration. *Ecol. Inform.* (2021).
- 583 49. Simmonds, M. B. *et al.* Guidelines for Publicly Archiving Terrestrial Model Data to Enhance
584 Usability, Intercomparison, and Synthesis. *Data Sci. J.* **21**, (2022).
- 585 50. Fowler, D., Barratt, J. & Walsh, P. Frictionless data: Making research data quality visible.

- 586 *Int. J. Digit. Curation* **12**, 274–285 (2018).
- 587 51. McNelis, J., Crow, M. & Devarakonda, R. ESS-DIVE File Level Metadata Extractor. *DOE*
588 *Code* <https://www.doi.org/10.11578/DC.20201103.5> (2020).
- 589 52. Jones, M. *et al.* Ecological Metadata Language version 2.2.0. *KNB Data Repository*
590 <https://www.doi.org/10.5063/F11834T2> (2019).
- 591 53. Sporny, M., Longley, D., Kellogg, G., Lanthaler, M. & Lindström, N. JSON-LD 1.0. *W3C*
592 *recommendation* **16**, 41 (2014).
- 593 54. Lehnert, K. A., Klump, J., Wyborn, L. & Ramdeen, S. IGSN: Trustworthy and Sustainable
594 Services for FAIR Samples. in *Geophysical Research Abstracts* vol. 21 (2019).
- 595 55. Yilmaz, P. *et al.* Minimum information about a marker gene sequence (MIMARKS) and
596 minimum information about any (x) sequence (MIxS) specifications. *Nat. Biotechnol.* **29**,
597 415–420 (2011).
- 598 56. CF Conventions and Metadata. <https://cfconventions.org/index.html> (2020).
- 599 57. Federal Geographic Data Committee. *Content Standard for Digital Geospatial Metadata*.
600 <https://ci.nii.ac.jp/naid/10016800076/> (1998).
- 601 58. Observations and Measurements. *Observations and Measurements*
602 <https://www.ogc.org/standards/om> (2011).
- 603 59. Walker, J. D., Lehnert, K. A., Hofmann, A. W., Sarbas, B. & Carlson, R. W. EarthChem:
604 International Collaboration for Solid Earth Geochemistry in Geoinformatics. in vol. 2005
605 IN44A-03 (2005).
- 606 60. Crystal-Ornelas, R. *et al.* Data from: “Enabling FAIR data in Earth and environmental
607 science with community-centric (meta)data reporting formats.” *ESS-DIVE*
608 <https://www.doi.org/10.15485/1866606> (2022).
- 609 61. JSON for Linking Data. *JSON for Linking Data* <https://json-ld.org/> (2016).
- 610 62. Brase, J. DataCite - A Global Registration Agency for Research Data. in *2009 Fourth*
611 *International Conference on Cooperation and Promotion of Information Resources in*

612 *Science and Technology* 257–261 (2009).

613 63. Evans, K. *et al.* *ASCII File Format Guidelines for Earth Science Data*. (2016).

614 64. Federal Geographic Data Committee. *National Geospatial Data Assets (NGDA) Metadata*
615 *Guidelines*. (2016).

616 65. Shafranovich, Y. *Common Format and MIME Type for Comma-Separated Values (CSV)*
617 *Files*. (2005).

618 66. USGS Data Dictionaries. <https://www.usgs.gov/data-management/data-dictionaries> (2020).

619 67. EDI. *Five phases of data publishing - Phase 2: Format and QC data tables*
620 <https://environmentaldatainitiative.org/five-phases-of-data-publishing/phase-2/> (2019).

621 68. Pepler, S. & Parton, G. BADC-CSV Format for Data Exchange.
622 <https://help.ceda.ac.uk/article/105-badc-csv> (2009).

623 69. Hsu, L. How “clean” should an Excel file be to be considered machine readable.
624 <https://my.usgs.gov/confluence/pages/viewpage.action?pageId=559852026> (2016).

625 70. NEON. NEON file naming conventions. [https://data.neonscience.org/file-naming-](https://data.neonscience.org/file-naming-conventions)
626 [conventions](https://data.neonscience.org/file-naming-conventions) (2022).

627 71. Tarboton, D. G., Horsburgh, J. S. & Maidment, D. R. CUAHSI community Observations
628 Data Model (ODM) version 1.1 design specifications. *Des Doc* (2008).

629 72. StreamPulse uploading data. http://pulseofstreams.weebly.com/uploading_data.html
630 (2022).

631 73. Cerf, V. *ASCII format for Network Interchange*. (1969).

632 74. Newell, D. B. & Tiesinga, E. *The International System of Units (SI) (2019 Edition)*. (National
633 Institute of Standards and Technology, 2019).

634 75. EPSG. *WGS 84*. (1984).

635 76. ORNL DAAC CSV Standards. <https://daac.ornl.gov/submit/csvstandards/> (2018).

636 77. Data Quality Review Checklist. <https://daac.ornl.gov/submit/qachecklist/>
637 <https://daac.ornl.gov/submit/qachecklist/> (2022).

- 638 78. USGS Data Templates. <https://www.usgs.gov/products/data-and-tools/data->
639 management/data-templates (2022).
- 640 79. National Archives. <https://www.archives.gov/records-mgmt/policy/transfer-guidance->
641 tables.html#structuredata (2022).
- 642 80. Klyne, G. & Newman, C. *Date and Time on the Internet: Timestamps*. (2002).
- 643 81. Loescher, H. AmeriFlux BASE CR-Lse La Selva. *AmeriFlux*
644 <https://www.doi.org/10.17190/AMF/1246013> (2016).
- 645 82. Torn, M. & Dengel, S. AmeriFlux US-NGB NGEE Barrow. *AmeriFlux*
646 <https://www.doi.org/10.17190/AMF/1436326> (2018).
- 647 83. Brokaw, N. Luquillo Forest Dynamics Plot (LFDP) Liana Data. *Environmental Data Initiative*
648 <http://www.doi.org/10.6073/PASTA/EABE6E15723324EA3938B456D5BB02C2> (2017).
- 649 84. Bueno de Mesquita, C. P. Plant colonization of moss-dominated soils in the alpine:
650 Microbial and biogeochemical implications. *Environmental Data Initiative*
651 <https://www.doi.org/10.6073/PASTA/C0CACD100CD89DA258B40A77FBB2FD4C> (2019).
- 652 85. Brooks, S. East Fork Poplar Creek sonde data at kilometer 16.2 water year 2019. *Mercury*
653 *Critical Interfaces SFA Data Search* <https://www.doi.org/10.12769/1569821> (2019).
- 654 86. Riscassi, A. & Brooks, S. East Fork Poplar Creek discharge at kilometer 5.4 water year
655 2012. *Mercury Critical Interfaces SFA Data Search* <https://www.doi.org/10.12769/1489524>
656 (2019).
- 657 87. National Ecological Observatory Network (NEON). Coarse downed wood bulk density
658 sampling (DP1.10014.001). *NEON* <https://www.doi.org/10.48443/Z0TG-QS14> (2020).
- 659 88. Salmon, V., Iversen, C., Childs, J. & VanderStel, H. NGEE arctic plant traits: Soil cores,
660 Kougatok road mile marker 64, Seward peninsula, Alaska, 2016. *NGEE Arctic Data Search*
661 <https://www.doi.org/10.5440/1346200> (2019).
- 662 89. Philben, M. *et al.* Results of experimental additions of organic nitrogen on soil organic
663 matter decomposition, teller road site, Seward peninsula, 2017 and 2018. *NGEE Arctic*

- 664 *Data Search* <https://www.doi.org/10.5440/1454263> (2019).
- 665 90. Yaffar, D., Lugo, A. E., Silver, W., Cuevas, E. & Molina Colon, S. Plant root trait
666 measurements raw data, 1962-2018, Island of Puerto Rico. *Next-Generation Ecosystem*
667 *Experiments Tropics; Oak Ridge National Laboratory*
668 <https://www.doi.org/10.15486/NGT/1558773> (2019).
- 669 91. Norby, R. *et al.* Root-soil depth profile in Luquillo Experimental Forest, Puerto Rico,
670 February, 2019. *Next-Generation Ecosystem Experiments Tropics; Oak Ridge National*
671 *Laboratory* <https://www.doi.org/10.15486/NGT/1574087> (2019).
- 672 92. Griffiths, N. & Sebestyen, S. SPRUCE porewater chemistry data for experimental plots
673 beginning in 2013. *Oak Ridge National Lab's Terrestrial Ecosystem Science Scientific*
674 *Focus Area (ORNL TES SFA)* <https://www.doi.org/10.3334/CDIAC/SPRUCE.028> (2016).
- 675 93. McPartland, M. Y. *et al.* SPRUCE: NDVI data from selected SPRUCE experimental plots,
676 2016-2018. *Oak Ridge National Laboratory (ORNL), Oak Ridge, TN (United States)*
677 <https://www.doi.org/10.25581/SPRUCE.057/1490190> (2019).
- 678 94. Croteau, M. N., Sikder, M., Poulin, B. A. & Baalousha, M. Laboratory data to assess the
679 effect of nanoparticle size and natural organic matter composition on the bioavailability of
680 platinum nanoparticles to a model freshwater invertebrate species. *U.S. Geological Survey*
681 <https://www.doi.org/10.5066/P9G18URX> (2020).
- 682 95. Danczak, R. E. *et al.* WHONDRS 48 Hour Diel Cycling Study at the Altamaha River in
683 Georgia, USA. *Worldwide Hydrobiogeochemistry Observation Network for Dynamic River*
684 *Systems (WHONDRS); ESS-DIVE* <https://www.doi.org/10.15485/1577263> (2019).
- 685 96. Stegen, J. C. *et al.* WHONDRS 48 hour Diel cycling study at HJ Andrews Experimental
686 Forest Watershed 1 (WS1). *Worldwide Hydrobiogeochemistry Observation Network for*
687 *Dynamic River Systems (WHONDRS); ESS-DIVE* <http://www.doi.org/10.15485/1509695>
688 (2019).
- 689 97. SESAR. SESAR Batch Registration Quick Guide. *Zenodo*

- 690 <https://www.doi.org/10.5281/zenodo.3874923> (2020).
- 691 98. IGSN Descriptive Metadata Schema. *IGSN metadata* <https://github.com/IGSN/metadata>
692 (2017).
- 693 99. Weibel, S., Kunze, J., Lagoze, C. & Wolf, M. Dublin core metadata for resource discovery.
694 *Internet Engineering Task Force RFC 2413*, 132 (1998).
- 695 100. ISO 19156:2011. *Geographic information — Observations and measurements*
696 <https://www.iso.org/standard/32574.html> (2011).
- 697 101. Joint Genome Institute Genome Online Database. <https://genome.jgi.doe.gov/portal/>
698 (2022).
- 699 102. USGS National Digital Catalog. [https://www.usgs.gov/programs/national-geological-and-](https://www.usgs.gov/programs/national-geological-and-geophysical-data-preservation-program/national-digital-catalog)
700 [geophysical-data-preservation-program/national-digital-catalog](https://www.usgs.gov/programs/national-geological-and-geophysical-data-preservation-program/national-digital-catalog) (2022).
- 701 103. Geologic Materials Repository Working Group. *The U.S. Geological Survey Geologic*
702 *Collections Management System (GCMS)—A master catalog and collections management*
703 *plan for U.S. Geological Survey geologic samples and sample collections.* (2015).
- 704 104. NEON Biorepository Data Portal. <https://biorepo.neonscience.org/portal/> (2022).
- 705 105. Hanson, B. Data policies and practices for AGU publications for models and model output.
706 (2020).
- 707 106. Williams, D. N., Lawrence, B. N., Lautenschlager, M., Middleton, D. & Balaji, V. The Earth
708 System Grid Federation: Delivering globally accessible petascale data for CMIP5.
709 *Proceedings of the Asia-Pacific Advanced Network* vol. 32 121 (2011).
- 710 107. Dryad. <https://datadryad.org/> (2022).
- 711 108. Zenodo. <https://zenodo.org/> (2022).
- 712 109. DAAC. <https://earthdata.nasa.gov/eosdis/daacs> (2021).
- 713 110. NCAR. <https://rda.ucar.edu/> (2022).
- 714 111. EOL. <https://www.eol.ucar.edu/about-eol> (2022).
- 715 112. Arctic Data Center. <https://arcticdata.io/> (2016).

- 716 113. Federal Geographic Data Committee. Content Standard for Digital Geospatial Metadata
717 Part 1: Biological Data Profile. (1999).
- 718 114. Christianson, D. S. *et al.* A metadata reporting framework (FRAMES) for synthesis of
719 ecohydrological observations. *Ecol. Inform.* **42**, 148–158 (2017).
- 720 115. USGS. <https://dashboard.waterdata.usgs.gov/> (2022).
- 721 116. SNOTEL. <https://wcc.sc.egov.usda.gov/reportGenerator/> (2022).
- 722 117. Earth Microbiome Project. [https://earthmicrobiome.org/protocols-and-standards/metadata-](https://earthmicrobiome.org/protocols-and-standards/metadata-guide/)
723 [guide/](https://earthmicrobiome.org/protocols-and-standards/metadata-guide/) (2022).
- 724 118. National Center for Biotechnology Information. SRA Metadata and Submission Overview.
725 <https://www.ncbi.nlm.nih.gov/sra/docs/submitmeta/> (2019).
- 726 119. Kattge, J. *et al.* TRY plant trait database - enhanced coverage and open access. *Glob.*
727 *Chang. Biol.* **26**, 119–188 (2020).
- 728 120. LeBauer, D. *et al.* BETYdb: a yield, trait, and ecosystem service database applied to
729 second generation bioenergy feedstock production. *Glob. Change Biol. Bioenergy* **10**, 61–
730 71 (2018).
- 731 121. Maitner, B. S., Boyle, B., Casler, N. & Condit, R. The bien r package: A tool to access the
732 Botanical Information and Ecology Network (BIEN) database. *Methods Ecol. Evol.* (2018).
- 733 122. ICOS. <https://www.icos-cp.eu/> (2022).
- 734 123. Pastorello, G. *et al.* The FLUXNET2015 dataset and the ONEFlux processing pipeline for
735 eddy covariance data. *Sci Data* **7**, 225 (2020).
- 736 124. Dorich, C. D. *et al.* Improving N₂O emission estimates with the global N₂O database. *Curr*
737 *Opin Environ Sustain* **47**, 13–20 (2020).
- 738 125. Lawrence, C. R. *et al.* An open-source database for the synthesis of soil radiocarbon data:
739 International Soil Radiocarbon Database (ISRaD) version 1.0. *Earth Syst. Sci. Data* **12**, 61–
740 76 (2020).
- 741 126. Schädel, C. *et al.* Decomposability of soil organic matter over time: the Soil Incubation

742 Database (SIDb, version 1.0) and guidance for incubation procedures. *Earth Syst. Sci. Data*
743 **12**, 1511–1524 (2020).

744 127. Jian, J. *et al.* A restructured and updated global soil respiration database (SRDB-V5). *Earth*
745 *Syst. Sci. Data* **13**, 255–267 (2021).

746 128. Ojima, D., Mosier, A., Del Grosso, S. & Parton, W. J. TRAGNET analysis and synthesis of
747 trace gas fluxes. *Global Biogeochem. Cycles* **14**, 995–997 (2000).

748 129. Hibbard, K. A., Law, B. E., Reichstein, M. & Sulzman, J. An analysis of soil respiration
749 across northern hemisphere temperate ecosystems. *Biogeochemistry* **73**, 29–70 (2005).

750 130. Horsburgh, J. S. *et al.* Observations Data Model 2: A community information model for
751 spatially discrete Earth observations. *Environ. Model. Soft.* **79**, 55–74 (2016).

752 131. NEON Protocols & Standardized Methods. [https://www.neonscience.org/data-](https://www.neonscience.org/data-collection/protocols-standardized-methods)
753 [collection/protocols-standardized-methods](https://www.neonscience.org/data-collection/protocols-standardized-methods) (2020).

754 132. EarthChem Data Templates. <https://www.earthchem.org/ecl/templates/> (2021).

755 133. SLAC. <https://www-ssrl.slac.stanford.edu/sfa/> (2022).

756 134. WFSFA. <https://eesa.lbl.gov/projects/watershed-function-sfa/> (2022).

757 135. River Corridor SFA. <https://www.pnnl.gov/projects/river-corridor> (2022).

758 136. AWH SFA. <https://www.anl.gov/bio/subsurface-biogeochemical-research> (2022).

759 137. LLNL Seaborg. <https://seaborg.llnl.gov/research/environmental-radiochemistry> (2022).

760 138. Mercury SFA. <https://www.esd.ornl.gov/programs/rsfa/> (2021).

761 139. WHONDORS. <https://www.pnnl.gov/projects/WHONDORS> (2022).

762 140. Ameriflux. <https://ameriflux.lbl.gov/> (2022).

763 141. NEON. <https://www.neonscience.org/> (2020).

764 142. Vogel, T. M. *et al.* TerraGenome: a consortium for the sequencing of a soil metagenome.
765 *Nat. Rev. Microbiol.* **7**, 252–252 (2009).

766 143. CUAHSI-HIS. *Master Controlled Vocabulary Registry for ODM 1.1*
767 <http://his.cuahsi.org/mastercvreg/cv11.aspx> (2008).

768 144. Blodgett, D., Lucido, J. & Kreft, J. Progress on water data integration and distribution: a
769 summary of select US Geological Survey data systems. *J. hydroinformatics* **18**, 226–237
770 (2016).

771 145. NWIS Inventory. <https://waterdata.usgs.gov/nwis/inventory> (2022).

772 146. NEON. [https://www.neonscience.org/data-samples/data-management/data-formats-](https://www.neonscience.org/data-samples/data-management/data-formats-conventions)
773 [conventions](https://www.neonscience.org/data-samples/data-management/data-formats-conventions) (2020).

774 147. WQP User Guide. https://www.waterqualitydata.us/portal_userguide/ (2022).

775 148. Garayburu-Caruso, V. A. *et al.* FTICR, NPOC, TN, and moisture of variably inundated
776 sediment across 48 north American rivers. *ESS-DIVE*
777 <https://www.doi.org/10.15485/1834208> (2021).

778 149. Alves, R. J. E. *et al.* Kinetic and temperature sensitivity properties of soil exoenzymes
779 through the soil profile down to one-meter depth at a temperate coniferous forest (Blodgett,
780 CA). *ESS-DIVE* <https://www.doi.org/10.15485/1830417> (2021).

781 150. Rogers, A., Ely, K. & Serbin, S. *Leaf Photosynthetic Parameters: Quantum Yield,*
782 *Convexity, Respiration, Gross CO₂ Assimilation Rate and Raw Gas Exchange Data,*
783 *Utqiagvik (Barrow), Alaska, 2016. NGEE Arctic Data Search*
784 <https://www.osti.gov/biblio/1482338> (2021).

785 151. Goldman, A. E., Emani, S. R., Pérez-Angel, L. C., Rodríguez-Ramos, J. A. & Stegen, J. C.
786 Perceived costs and benefits of ICON science and foundational documents associated with
787 “Integrated, Coordinated, Open, and Networked (ICON) science to Advance the
788 Geosciences: Introduction and Synthesis of a Special Collection of Commentary Articles”.
789 (2022).

790 152. Roebuck, A. *et al.* FTICR-MS, Sensor, and Environmental Data from 5 Streams Impacted
791 by the 2020 Holiday Farm Fire Associated with: “Spatiotemporal controls on the delivery of
792 dissolved organic matter to streams following a wildfire.” *ESS-DIVE*
793 <https://www.doi.org/10.15485/1869708> (2022).

- 794 153. Allison, S. & Martiny, J. B. H. Fungal and bacterial growth variation due to drought and
795 nitrogen addition experimental treatments. Loma Ridge Experimental Project. 2010-2012.
796 *ESS-DIVE* <https://www.doi.org/10.15485/1828589> (2021).
- 797 154. Dove, N., Torn, M., Hart, S. & Tas, N. Chemistry data from soils and soil incubation
798 experiments from the whole-soil warming experiment at Blodgett Forest, CA, 2018, from:
799 “Metabolic capabilities mute positive response to direct and indirect impacts of warming
800 throughout the soil profile.” *ESS-DIVE* <https://www.doi.org/10.15485/1866269> (2022).
- 801 155. ESS-DIVE. *ESS-DIVE* <https://ess-dive.lbl.gov/> (2022).
- 802

Tables and Figures

Table 1

A list of all 11 ESS-DIVE (meta)data reporting formats and citations for associated publications. Each row includes a description of the reporting format, a list of existing data standards as well as other resources (i.e., data repositories, data systems, datasets, data files from DOE projects, instruments) consulted during development of each reporting format.

| Reporting format | Description | Standards reviewed | Other resources reviewed |
|--------------------------------------|--|--|---|
| Dataset metadata ²⁵ | The dataset metadata reporting format is in use as part of ESS-DIVE's data submission process. It ensures that data submitted to ESS-DIVE has a minimum set of metadata (e.g., title, abstract, author contact information) to enable data reuse. | <ul style="list-style-type: none"> • Ecological metadata language⁵² • JavaScript Object Notation for Linked Data⁶¹ • DataCite 4.1⁶² • Federal Geographic Data Committee⁵⁷ | <ul style="list-style-type: none"> • Office of Science and Technical Information Announcement Notice 241.6⁴³ |
| File-level metadata ^{26,51} | Metadata about the data files included as part of a larger dataset. Information about the data files and the spatiotemporal context for each data file. Includes script that can extract file-level metadata from data files that adhere to this format. | <ul style="list-style-type: none"> • Ecological Metadata Language⁵² • ASCII File Format Guidelines for Earth Science Data⁶³ • National Geospatial Data Assets Metadata Guidelines⁶⁴ • Common format for CSV files⁶⁵ • USGS Data Dictionaries⁶⁶ | <ul style="list-style-type: none"> • Environmental Data Initiative: Five phases of data publishing⁶⁷ • Centre for Environmental Data Analysis⁶⁸ • United States Geological Survey: How "clean" should Excel files be to be considered machine readable⁶⁹ • National Ecological Observatory Network file naming conventions⁷⁰ • Consortium of Universities for the Advancement of Hydrologic Science ODM specifications v1.1.1⁷¹ • StreamPulse⁷² |

| | | | |
|---|--|---|---|
| | | | <ul style="list-style-type: none"> • Ameriflux³⁹ |
| Comma Separated Value files ²⁷ | Recommendations for organizing tabular data stored in Comma Separated Values. Recommends optimal file structure, naming structure, and field structure. | <ul style="list-style-type: none"> • Common format for CSV files⁶⁵ • ASCII format for Network Interchange⁷³ • International System of Units⁷⁴ • World Geodetic System 1984⁷⁵ • ORNL DAAC CSV Standards⁷⁶ • ORNL DAAC Data Quality Checklist⁷⁷ • USGS Data Templates⁷⁸ • National Archives: Tables of File Formats⁷⁹ • ASCII File Format Guidelines for Earth Science Data⁶³ • Date and Time on the Internet: Timestamps⁸⁰ | <p><i>Datasets downloaded from repositories/databases:</i></p> <ul style="list-style-type: none"> • Ameriflux^{81,82} • Environmental data initiative^{83,84} • Mercury SFA^{85,86} • NEON data product⁸⁷ • Next-Generation Ecosystem Experiments - Arctic^{88,89} • Next Generation Ecosystem Experiments - Tropics^{90,91} • Spruce and Peatland Responses Under Changing Environments^{92,93} • United States Geological Survey⁹⁴ • WHONDRS^{95,96} <p><i>Other data resources:</i></p> <ul style="list-style-type: none"> • National Ecological Observatory Network file naming conventions⁷⁰ • Environmental Data Initiative: Five phases of data publishing⁶⁷ • Centre for Environmental Data Analysis⁶⁸ |
| Sample ID ^{28,44} and metadata | Enables tracking and reuse of environmental samples by providing guidance for registering physical samples for persistent identifiers. Enables sample linking across multiple collaborators, labs, and data systems online. Includes sample identifier metadata, descriptions about the sample and location. | <ul style="list-style-type: none"> • SESAR IGSN⁹⁷ • IGSN Schema⁹⁸ • DataCite⁶² • Dublin Core⁹⁹ • Darwin Core⁸ • Minimum information about any sequence (MIxS)⁵⁵ • Observations and measurements¹⁰⁰ | <ul style="list-style-type: none"> • Joint Genome Institute Genome Online Database and sample template¹⁰¹ • Environmental Protection Agency Water Quality Exchange¹⁴ • United States Geological Survey National Digital Catalog¹⁰² • National Geological and Geophysical Data Preservation Program Geologic Collections |

| | | | |
|---|---|--|---|
| | | | <p>Management System¹⁰³</p> <ul style="list-style-type: none"> • National Ecological Observatory Network Biorepository Data Portal¹⁰⁴ |
| <p>Model archiving guidelines^{29,49}</p> | <p>Helps modelers determine which components of their model data to archive and how to bundle data files for publication based on authorship criteria and repository storage capacities. Recommendations include archiving model inputs, testing or validation data, and any code or scripts associated with data preprocessing or modeling runs.</p> | <ul style="list-style-type: none"> • National Science Foundation Earthcube modeling rubric¹⁰⁵ | <ul style="list-style-type: none"> • Earth System Grid Federation¹⁰⁶ • Dryad¹⁰⁷ • Zenodo¹⁰⁸ • NASA's Earth Observing System Data and Information System¹⁰⁹ • National Center for Atmospheric Research Data Archive¹¹⁰ • Earth Observatory Lab data archive¹¹¹ • ESS-DIVE⁴ • National Science Foundation Arctic Data Center¹¹² |
| <p>Locations metadata³⁰</p> | <p>Minimum set of metadata to report about locations where field samples or observations are collected. While the sampling features described in the reporting format are primarily point locations, the format can optionally be expanded to capture metadata about plots and regions</p> | <ul style="list-style-type: none"> • CF Conventions⁵⁶ • DarwinCore⁸ • Content Standard for Digital Geospatial Metadata⁵⁷ • Content Standard for Digital Geospatial Metadata: Biological Data Profile¹¹³ • Open Geospatial Consortium observations and measurements⁵⁸ | <ul style="list-style-type: none"> • Watershed Function Scientific Focus Area sensor metadata template⁴¹ • Watershed function SFA locations template⁴² • Ameriflux Biological, Ancillary, Disturbance, and Metadata¹⁵ • iSamples⁴⁶ • Framework for Reporting Data and Metadata for Earth Science¹¹⁴ • United States Geological Survey National Water Dashboard¹¹⁵ • United States Department of Agriculture Snowpack Telemetry¹¹⁶ |
| <p>Amplicon abundance</p> | <p>Microbial abundance profiles and sequences obtained from genomic</p> | <ul style="list-style-type: none"> • Minimum information about a marker gene sequence | <ul style="list-style-type: none"> • Earth Microbiome Project Metadata Guide¹¹⁷ |

| | | | |
|--|---|---|--|
| table metadata ³¹ | sequencing of environmental samples | (MIMARKS) for sequencing metadata. ⁵⁵ <ul style="list-style-type: none"> No data standards existed for amplicon abundance table bioinformatic processing metadata. | <ul style="list-style-type: none"> National Center for Biotechnology Information (NCBI) Sequence Read Archive (SRA) Metadata and Submission Overview¹¹⁸ |
| Leaf-level gas exchange ^{32,47} | Guidelines for providing methods metadata and instrument metadata for experiments collecting leaf-level gas exchange data. Defined variable names and units for data, and specifies required variables for various measurement types. | No data standards existed | <ul style="list-style-type: none"> TRY database¹¹⁹ Biofuel Ecophysiological Traits and Yields¹²⁰ Botanical Information and Ecology Network¹²¹ <p><i>Instruments</i></p> <ul style="list-style-type: none"> ADC iFL ADC LCi T ADC proT CID CI-340 Licor 6400XT Licor 6800 PP Systems CIRAS-2 PP Systems CIRAS-3 PP Systems TARGAS-1 Walz GFS-3000 |
| Soil respiration ^{33,48} | Brings together decentralized and individual-driven field-measured data. The chamber-level metadata template facilitates reporting information about experimental conditions including variables being measured, chamber size and design, and instrumentation. The other template is for data entry of CO ₂ measurements from experimental chambers. | <ul style="list-style-type: none"> Integrated Carbon Observation System¹²² | <ul style="list-style-type: none"> AmeriFlux³⁸ FLUXNET¹²³ Global N₂O Database¹²⁴ International Soil Radiocarbon Database¹²⁵ Soil Incubator Database¹²⁶ Soil Respiration Database¹²⁷ Trace Gas Network¹²⁸ Hibbard et al. 2005¹²⁹ |

| | | | |
|---|---|---|---|
| <p>Sample-based water and soil chemistry³⁴</p> | <p>Recommendations for consistently formatting data and metadata obtained from laboratory chemical analysis of field samples. Reporting chemical concentrations and metadata for data collection methods.</p> | <ul style="list-style-type: none"> ● Observations Data Model 2 (ODM2)¹³⁰ ● Minimum information about any sequence⁵⁵ ● Water Quality Exchange¹⁴ ● National Ecological Observatory Network Protocols & Standardized Methods¹³¹ ● EarthChem¹³² | <ul style="list-style-type: none"> ● SLAC Floodplain Hydro-Biogeochemistry SFA¹³³ ● LBNL Watershed Function SFA¹³⁴ ● PNNL River Corridor Hydrobiogeochemistry SFA¹³⁵ ● ANL Wetlands SFA¹³⁶ ● LLNL Biogeochemistry at Interfaces SFA¹³⁷ ● ORNL Mercury SFA¹³⁸ ● WHONDRS¹³⁹ ● AmeriFlux¹⁴⁰ ● NEON¹⁴¹ ● Terragenome¹⁴² |
| <p>Water level and sonde-based hydrologic monitoring³⁵</p> | <p>Hydrologic sensor-based data and corresponding deployment metadata. Includes a terminology file that recommends terms for commonly used hydrological variable names that had significant variation across prior standards and data resources</p> | <ul style="list-style-type: none"> ● Observations Data Model 1 used in the Environmental Data Initiative (EDI) and Consortium of Universities for the Advancement of Hydrologic Science - Hydrologic Information System¹⁴³ ● Observations Data Model 2¹³⁰ | <ul style="list-style-type: none"> ● United States Geological Survey National Water Information Systems^{144,145} ● National Ecological Observatory Network Data Formats and Conventions¹⁴⁶ ● Water Quality Portal¹⁴⁷ |

Table 2

Examples of datasets published on ESS-DIVE utilizing at least one of the 11 ESS-DIVE (meta)data reporting formats. Each row includes the dataset title, citation, and the reporting format(s) used in the dataset.

| Dataset Title | Reporting Format(s) Used |
|--|---|
| FTICR, NPOC, TN, and Moisture of Variably Inundated Sediment across 48 North American Rivers ¹⁴⁸ | Sample-based water and soil chemistry, Sample ID and metadata, Comma Separated Value files, and File-level metadata Reporting Formats |
| Kinetic and temperature sensitivity properties of soil exoenzymes through the soil profile down to one-meter depth at a temperate coniferous forest (Blodgett, CA) ¹⁴⁹ | Sample ID and metadata, Comma Separated Value files, and File-level metadata Reporting Formats |
| Leaf Photosynthetic Parameters: Quantum Yield, Convexity, Respiration, Gross CO ₂ Assimilation Rate and Raw Gas Exchange Data, Utqiagvik (Barrow), Alaska, 2016 ¹⁵⁰ | Leaf-level gas exchange Reporting Format |
| Perceived Costs and Benefits of ICON Science and Foundational Documents associated with "Integrated, Coordinated, Open, and Networked (ICON) Science to Advance the Geosciences: Introduction and Synthesis of a Special Collection of Commentary Articles" ¹⁵¹ | File-level metadata Reporting Format |
| FTICR-MS, Sensor, and Environmental Data from 5 Streams Impacted by the 2020 Holiday Farm Fire Associated with: "Spatiotemporal controls on the delivery of dissolved organic matter to streams following a wildfire" ¹⁵² | Hydrologic Monitoring, Comma Separated Value files, and File-level metadata Reporting Formats |
| Fungal and bacterial growth variation due to drought and nitrogen addition experimental treatments. ¹⁵³ | File-level metadata and Comma Separated Value files Reporting Formats |
| Chemistry data from soils and soil incubation experiments from the whole-soil warming experiment at Blodgett Forest, CA, 2018, from: "Metabolic capabilities mute positive response to direct and indirect impacts of warming throughout the soil profile" ¹⁵⁴ | File-level metadata and Comma Separated Value files Reporting Formats |

Table 3

Evaluation of how ESS-DIVE’s reporting formats achieve the Findable, Accessible, Interoperable, and Reusable data principles ^{1,21} in the context of the ESS-DIVE data repository (some of the FAIR Principles are not applicable at the reporting format level and implementation occurs through the data repository supported infrastructure / metadata). All datasets published in the ESS-DIVE repository adhere to the Dataset Metadata reporting format, which at least partially satisfies many FAIR principles, particularly Findability and Accessibility. For some FAIR principles (e.g., F4) submitting data to ESS-DIVE (regardless of which domain-specific format is used) will automatically ensure that the principle is met. When reporting formats do not fully achieve a FAIR principle, we list the reporting format and progress made toward fully machine-actionable and FAIR (meta)data.

| FAIR Principle | Number of reporting formats achieving FAIR principle when combined with repository features | Reporting formats not fully achieving FAIR principle | Progress toward achieving FAIR principle |
|--|---|---|---|
| F1. (meta)data are assigned a globally unique and persistent identifier. | 11 | All achieve FAIR principle on ESS-DIVE or other data repositories | All datasets uploaded to ESS-DIVE are associated with a persistent DOI. Data and metadata associated with samples and associated metadata are assigned an IGSN. |
| F2. data are described with rich metadata | 11 | All achieve FAIR principle | All reporting formats provide guidance and many include templates for including descriptive metadata about entire datasets, individual data files, samples, and specific data types. Metadata requirements were defined by scientists who provide and potentially search for the data. |
| F3. (meta)data clearly and explicitly include the identifier of the data it describes. | 9 | CSV reporting format*, model data archiving guidelines* | Data files are assigned an identifier and explicitly linked via ESS-DIVE’s use of the Ecological Metadata Language ⁹² |
| F4. (meta)data are registered or indexed in a searchable resource. | 1 | CSV reporting format, model data archiving guidelines, soil respiration, leaf-level gas exchange, sonde-based hydrologic monitoring, sample-based water and soil chemistry, amplicon abundance table metadata, file-level | Metadata provided about each dataset submitted to ESS-DIVE are indexed and searchable as part of the ESS-DIVE repository. Individual reporting formats are not currently at a stage where the data contained within files are automatically indexed and searchable except for dataset metadata. This is a target for future work. |

| | | | |
|--|----|--|---|
| | | metadata, sample ID and metadata, locations metadata | |
| A1. (meta)data are retrievable by their identifier using a standardized communications protocol. | 9 | CSV reporting format*, model data archiving guidelines* | When data are archived in the ESS-DIVE repository, metadata are available over HTTP. |
| A1.1. the protocol is open, free, and universally implementable. | 11 | All achieve FAIR principle on ESS-DIVE or other data repositories | Users can access (meta)data on ESS-DIVE over HTTP. |
| A1.2. the protocol allows for an authentication and authorization procedure, where necessary. | 11 | All achieve FAIR principle on ESS-DIVE | Data that adheres to reporting formats is stored in ESS-DIVE repository which allows for authentication and authorization, so individual reporting formats do not need to have their own protocols. |
| A2. (meta)data are accessible, even when the data are no longer available | 11 | All achieve FAIR principle on ESS-DIVE | Data are stored using persistent identifiers and in non-proprietary formats. Also, ESS-DIVE stores metadata independent of the datasets and the metadata are replicated across the DataONE network. |
| I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation. | 1 | CSV reporting format, model data archiving guidelines, soil respiration, leaf-level gas exchange, sonde-based hydrologic monitoring, sample-based water and soil chemistry, amplicon abundance table metadata, file-level metadata, sample ID and metadata, locations metadata | All reporting formats made good strides toward machine readable knowledge representation. However, based on community feedback most prioritized human-readable explanations of variables via data dictionaries and term lists rather than fully machine readable knowledge representation. Having flexibility in term usage was considered important to lower the barrier for data contributors to adopt the reporting formats. |
| I2. (meta)data use vocabularies that follow FAIR principles. | 4 | CSV reporting format, model data archiving guidelines, soil respiration, leaf-level gas exchange, amplicon abundance table metadata, file-level metadata, locations metadata | Based on feedback from our community engagement, we made pragmatic choices that involved leaving out ontologies that made adoption of the formats burdensome in lieu of terms that scientists preferred to use. For samples, we chose a small subset of formal ontology terms to describe environmental context relevant for the ESS community. Scientists |

| | | | |
|---|----|----------------------------|--|
| | | | are encouraged, but not required to use ontologies when applicable. |
| I3. (meta)data include qualified references to other (meta)data. | 11 | All achieve FAIR principle | All data submitted to ESS-DIVE can be externally linked to related resources in the repository. |
| R1. (meta)data are richly described with a plurality of accurate and relevant attributes. | 11 | All achieve FAIR principle | The dataset metadata and file-level metadata reporting formats ensure that entire datasets and individual files are described with rich metadata. Moreover, reporting formats also provide space for metadata provision for specific data types (e.g., chamber metadata for the soil respiration reporting format). |
| R1.1. (meta)data are released with a clear and accessible data usage license. | 11 | All achieve FAIR principle | All (meta)data submitted to ESS-DIVE are published with either CC BY 0 or CC BY 4.0 licenses. |
| R1.2. (meta)data are associated with detailed provenance [^] . | 11 | All achieve FAIR principle | All reporting formats make progress at tracking information about how and/or when data were collected via elements that collect detailed information about locations and timestamps when data were collected. Some reporting formats (e.g., sonde-based hydrologic monitoring, sample-based water and soil chemistry) include templates for methods-specific metadata. ESS-DIVE allows linking datasets to associated external (meta)data using a set of formal machine-readable relationship fields/terms, and is working towards extending to a variety of related resources involved in ESS work. |
| R1.3. (meta)data meet domain-relevant community standards | 11 | All achieve FAIR principle | The community reporting format effort is directly achieving this principle by involving research communities in creating (meta)data standards (i.e., reporting formats). |

An asterisk (*) indicates when the FAIR principle did not apply to the reporting format because it was an instruction-based guideline rather than templates and tools.

The caret (^) indicates that our interpretation of provenance in this case is that all reporting formats allow users to track components of provenance fundamentals including timestamps, location information, and instruments used, among other elements.

A plus (+) indicates that the FAIR principle is met when combined with specific features offered by the ESS-DIVE repository as specified in the last column of the table. The reporting format can achieve the FAIR principle on other repositories that offer the same features (e.g., assignment of persistent identifiers).

Box 1

Guidelines for research communities to self-organize and create, document, and share (meta)data reporting formats when formats do not already exist or do not fit scientists' needs.

1. **Research existing (meta)data standards and other data resources** across agencies and organizations both within the US and internationally.
2. Create a **(meta)data crosswalk** (Supplementary Files 1-10) to define how other standards and data resources translate to the proposed reporting format.
3. Work with the scientific community to **iteratively develop and obtain feedback** (see Figure 2) on (meta)data reporting format.
4. Develop **documentation** (instructions, templates, variables, descriptions, units, metadata) to support the format. Consider appropriate **file formats** for any templates.
5. Archive finalized version of the reporting format in a **long-term data repository** as well as a version control platform (e.g., GitHub³⁷).

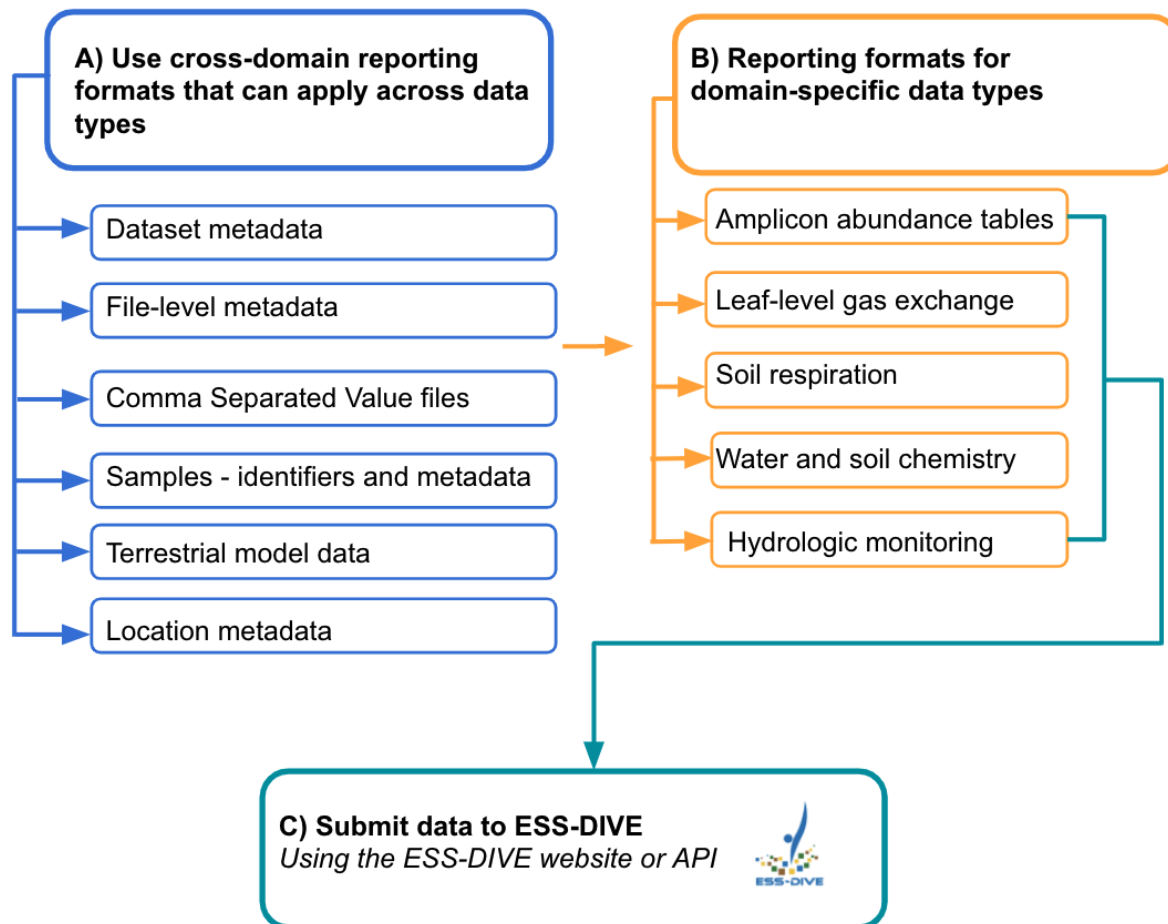


Figure 1

Workflow to help determine which (meta)data reporting formats apply to datasets. The set of 11 ESS-DIVE (meta)data formats are either (A) cross-domain guidelines that can be applied to many data types or (B) are data type-specific. For those archiving data with ESS-DIVE, researchers can upload data through the ESS-DIVE web user interface¹⁵⁵ or programmatically through an API.



Figure 2

Each of the 11 ESS-DIVE (meta)data reporting formats were developed in cross-functional teams that often involved domain scientists, software engineers, and informatics specialists.