1    Implementation of a workflow for publishing citeable environmental data: successes,

2    challenges and opportunities from a data centre perspective.

3

4    Kathryn A. Harrison, Daniel G. Wright, Philip Trembath

5

6    **Affiliation**: Centre for Ecology & Hydrology, Lancaster Environmental Centre, Library

7    Avenue, Bailrigg, Lancaster, LA1 4AP, UK

8

9    **Corresponding author**: Kathryn A. Harrison

10   **Email:** kath@ceh.ac.uk

11

12

13

14    Abstract

15    In recent years, the development and implementation of a robust way to cite data has

16    encouraged many previously sceptical environmental researchers to publish the data they

17    create, thus ensuring more data than ever are now open and available for re-use within and

18    between research communities. Here we describe a workflow for publishing citeable data in

19    the context of the environmental sciences – an area spanning many domains and generating

20    a vast array of heterogeneous data products. The processes and tools we have developed

21    have enabled rapid publication of quality data products including datasets, models and

22    model outputs which can be accessed, re-used and subsequently cited. However, there are

23    still many challenges that need to be addressed before researchers in the environmental

24    sciences fully accept the notion that datasets are valued outputs and time should be spent in

25    properly describing, storing and citing them. Here we identify current challenges such as

26    citation of dynamic datasets and issues of recording and presenting citation metrics. In

27    conclusion, whilst data centres may have the infrastructure, tools, resources and processes

28    available to publish citeable datasets, further work is required before large-scale uptake of

29    the services offered is achieved. We believe that once current challenges are met, data

30    resources will be viewed similarly to journal publications, as valued outputs in a researcher's

31    portfolio, and therefore both the quality and quantity of data published will increase.

32

33

36

37      1.0 Introduction

38      Historically, there has been resistance from some researchers in the environmental sciences

39      to publishing data, other than referring to it in articles in recognised scientific journals. The

40      act of making data openly available for the public to view, access and re-use is an unfamiliar

41      concept to many, although, for some scientific communities (e.g. bioinformatics and 'omics)

42      data archival is a cultural norm [1]. Inability to access scientific data is an obstacle to

43      interdisciplinary research [2, 3] which is key in the area of the environmental sciences as

44      they cover a broad range of disciplines and often aim to answer complex questions requiring

45      input from a range of specialists. Whilst each year large amounts of research funds

46      (including tax payers' money) are spent generating new data, existing data remain

47      inaccessible, unidentified and therefore underutilised [4].

48

49      In recent years there has been increasing pressure on scientists to make the data they

50      generate openly available. Regulatory pressures such as the EU's INSPIRE[1] directive and

51      compliance with research funders' policies (e.g. RCUK'[2] data policy) are compelling

52      researchers to publish their data. Nonetheless, this regulatory approach has done little to

53      prompt a significant change in cultural practices. It is clear that in order for a shift in

54      behaviour to occur, researchers must feel confident that making the data they create

55      available will not adversely impact on their career. Mayernik [5] and Assante et al [6] make

56      reference to the cultural barriers which make scientists unwilling to share results and

57      document the fears that researchers have of being 'scooped', their data being used

58      improperly and the difficulties they face in producing data in a shareable form.

59

---

[1] http://inspire.ec.europa.eu/
[2] http://www.rcuk.ac.uk/research/datapolicy/

60    If a published data resource is regarded as a citeable publication it can impact positively on

61    a researcher's reputation [2] and this in turn will encourage the publication of more data.

62    Generation and subsequent publication of data should be recognised as valuable activities

63    but currently lacks an essential pre-requisite – accepted metrics of significance [7]. For

64    example, it should be possible to collect information on who has re-used the data, what it

65    has been re-used for and how many times has it been re-used. Metrics such as these could

66    ultimately bear on the academic reputation of a researcher amongst their peers in a similar

67    way that metrics on citations of journal papers currently do. Provision of this service alone

68    will not solve all the problems, however, and it will take time to establish. Data centres and

69    research institutions must also consider providing support to researchers, increasing

70    awareness of the issues and developing simple workflows in order that time-limited

71    researchers can engage in the process of making the data they create publicly available and

72    gain credit for doing so.

73

74    Providing a means of citing data allows data creators to be perpetually linked to the datasets

75    they produce. However, for researchers to gain credit for their work a formal, community-

76    recognised structure must be set in place [2]. DataCite[3] has been instrumental in developing

77    and supporting the standards behind persistent identifiers for data. They provide a means by

78    which researchers can find, identify and cite research data and other research objects.

79    DataCite currently use the Digital Object Identifier (DOI) system[4] as a persistent identifier for

80    data resources, although other permanent identifiers could be used in a similar way

81    [8].Through this system, DataCite is able to provide a robust mechanism for allowing citation

82    of data resources. The DOI system is one of the more suitable candidates for permanently

83    identifying research data the system is well-established and widely used for identifying

84    research articles and are therefore a familiar entity to researchers [9, 10] and publishers

---

[3] https://www.datacite.org/
[4] http://www.doi.org/

85     alike. Whilst a suitable system for identifying data and making it citeable is in place there still

86     exists a gap between the technical ability to cite data and the cultural behaviour of

87     researchers within the environmental sciences (see above). This gap can only be narrowed

88     by researchers interacting with the system in a positive manner and gaining reward for doing

89     so, for example, a raised awareness of a researcher's work within the community leading to

90     increased collaborative or funding opportunities or improved promotion prospects [3].

91

92     The Environmental Information Data Centre (EIDC[5]) is a Natural Environment Research

93     Council (NERC) environmental data centre hosted by the Centre for Ecology and Hydrology

94     (CEH). The data centre primarily accepts data resources from NERC-funded research

95     covering a wide spectrum of disciplines including the terrestrial and freshwater sciences and

96     hydrology. Data held by the EIDC is usually 'complete' end of project life data, although the

97     data centre also holds data collected from long-term environmental monitoring programmes

98     – normally deposited in discrete time slices. The EIDC offers researchers the opportunity to

99     obtain a DOI for data they have created and therefore the ability to cite the resource in

100    literature. DOIs are used as a permanent identifier for data held by the EIDC as this is the

101    identifier initially chosen by NERC for use in its data centres. NERC works with The British

102    Library who is the allocation agent for DataCite in the UK. By assigning a DOI to a resource,

103    the EIDC are signifying that the data are complete, stable, in a useable format, have

104    appropriate metadata, have passed the quality control checks within the domain expertise of

105    the data centre and have guaranteed long-term curation [9]. Whilst there is nothing inherent

106    in a DOI that guarantees the data it identifies will remain permanently available and stable,

107    the EIDC holds a form of 'social contract' between itself and the registry (DataCite and the

108    British Library)  to ensure that this is the case [10, 11]. The EIDC uses checksums to ensure

109    data remain unchanged once they have been deposited with the data centre and data

---

[5] http://eidc.ceh.ac.uk/

110    depositors receive a copy of the checksum so that they may verify this at any given time. As

111    a data centre, the EIDC has been offering DOIs for datasets that it holds since 2011. Here

112    we outline the processes established to provide this service and describe initial community

113    use and acceptance of the system. We explore the impact that this service has had on the

114    data centre, the datasets published by the data centre and the subsequent exposure of

115    those datasets. Further, we discuss future challenges for the EIDC, specifically, citation of

116    dynamic datasets and the collection of citation metrics. Both these issues have the ability to

117    further influence the volume and quality of data published within the environmental sciences

118    community.

119

120    2.0 Data centre process for obtaining a DOI

121    Data resources are taken into the EIDC following a defined workflow, which includes strict

122    process and quality control measures. Data resources which are identified as suitable for

123    deposit are curated by the data centre in order that they may be viewed and accessed over

124    the long-term. For a data resource to be deemed suitable it must meet a number of criteria

125    such as subject area, funder, repeatability and uniqueness – data held elsewhere would not

126    be considered for deposit. The EIDC first began using a defined workflow in 2011 and to-

127    date holds a total of nearly 400 data resources including datasets, models, model outputs

128    and web services. Only datasets that have passed through the workflow and been formally

129    'ingested' into the EIDC are eligible for a DOI. Each of the seven NERC data centres (of

130    which EIDC is one) has a representative who can register DOIs for NERC datasets. Whilst

131    the act of registering a DOI with DataCite is the same for all data centres, the manner in

132    which datasets are prepared to a form which is acceptable for allocation of a DOI varies.

133

134    2.1. Support for researchers

135     The EIDC is hosted by the Centre for Ecology and Hydrology (CEH) and as such, the data

136     centre accepts data from both 'internal' depositors (i.e. researchers from CEH) and 'external'

137     depositors (i.e. researchers employed elsewhere such as universities and other research

138     institutes). The process for ingesting data is identical for both internal and external

139     depositors; the support given to researchers prior to submission of the data is also broadly

140     similar and will be described here briefly. CEH employ a team of Informatics Liaison (IL) staff

141     whose role it is to support researchers with data management and all that it entails.

142     Members of the IL team will work with researchers ideally from the very start of a project to

143     ensure a data management plan is created and regularly reviewed and updated. Likewise,

144     this support is also available for 'external' researchers whose data will ultimately be

145     considered for deposit with the EIDC. Data management plans identify the data resources

146     that will be offered to EIDC and also list the supporting documentation which will accompany

147     the deposit. The IL staff will initiate a deposit once the researcher is ready, and support them

148     through the process – for example, helping to complete discovery metadata records, giving

149     advice on formatting the data for deposit, creating any supporting documentation and

150     discussing issues such as licensing and citation. The workflow whereby the EIDC registers

151     DOIs for data it holds is described below. However, this workflow does not solely include

152     actions carried out automatically by the data centre with the researcher in isolation, support

153     from IL and data centre staff is provided throughout. A full description of the complete

154     workflow for ingesting data into the data centre is not within the scope of this article but has

155     been described elsewhere within this special issue.

156

157     2.2. Discovery metadata

158     At the EIDC, the process for obtaining a DOI begins with the collection and storage of

159     discovery metadata. The EIDC uses the UK GEMINI[6] metadata specification for describing

---

[6] http://www.agi.org.uk/join-us/agi-groups/standards-committee/uk-gemini

160   the data resources for discovery purposes. This standard has a set of mandatory

161   requirements and includes elements such as title, abstract, lineage and keywords. The

162   EIDCs discovery-level metadata is stored in a metadata file store based on Git[7], a distributed

163   revision control system, which ensures a complete history of all changes made to metadata

164   is maintained. Metadata are stored as JSON[8], an open data-interchange format that records

165   data as attribute-value pairs. The JSON format allows the data centre to transform the data

166   and present them in a number of different formats targeted at distinct audiences – being both

167   human- and machine- readable. For example, the metadata can be presented as a human

168   readable HTML page, as GEMINI-compliant XML for data exchange to data.gov or as XML

169   in the DataCite schema[9] for registering DOIs and populating the DataCite catalogue.

170   Metadata records are created by the researcher depositing data with help from data centre

171   staff, who enter the information using a bespoke metadata editing tool; the metadata is

172   accessed from the CEH catalogue[10]. This catalogue was developed in-house to provide the

173   public with a user-friendly interface for finding, viewing and accessing data (Fig 1).

174

175   The discovery metadata record for the data resource also acts as the landing page for the

176   DOI, once registered, and it was designed with this in mind. Although much of the

177   information about a resource is captured using the GEMINI metadata standard, how it

178   should be presented to function as a DOI landing page was key to the decisions made about

179   the way the page was fashioned. As stated by Ball and Duke [12] a landing page should

180   "enable readers to ensure they have located the right dataset, to (re-)familiarize themselves

181   with the research context and supporting documentation, to consider licence terms prior to

182   downloading and to switch to a more recent version of the data if required" (pg 12). The

183   EIDC is keen to promote the use of data citations, therefore, once a resource has a DOI, this

---

[7] http://www.git-scm.com/
[8] http://www.json.org/
[9] http://schema.datacite.org/
[10] https://catalogue.ceh.ac.uk/

184    appears, together with the reference to be quoted with any subsequent re-use, at the very

185    top of the page, immediately below the title. An abstract describing the resource follows the

186    DOI and to the right, in a 'Get the data' panel, information on how to order the data, access

187    to supporting documentation and another full citation for the data is presented with the clear

188    instruction 'If you reuse this data you must cite:' (Fig. 1). In designing the landing page,

189    particular care was taken to use accessible language rather than adopt the somewhat

190    opaque language of the metadata standard. For example, 'resource locator' is labelled

191    'online resources' and 'responsible organisation' is labelled 'contacts'. The GEMINI XML

192    view of the metadata retains the standard terms, it is solely the landing page/catalogue view

193    that presents the more user-friendly version.

194

195    2.3. DOI registration

196    The CEH Catalogue generates DataCite metadata directly from the GEMINI metadata using

197    a simple mapping (Table 1). To register a new DOI, the designated DOI administrator makes

198    a request by clicking a hyperlink on the data resource's record in the data catalogue. This

199    hyperlink only appears on the record if a number of conditions are met. First, only a DOI

200    administrator has access to the link – it does not appear if a user without the necessary

201    permissions is logged in. Second, all the key pre-requisite elements of DataCite metadata

202    must be present within the record – namely: at least one author; a date of publication; a title;

203    and a publisher (other information is also included in the DataCite metadata but these are

204    the only mandatory fields). Thirdly, the landing page must be publicly accessible. Fourth and

205    finally, there must not already be a DOI registered for that resource.  By clicking the

206    hyperlink, this triggers a series of actions which occur programmatically without the need for

207    further user intervention. The metadata is posted to DataCite's REST API[11], this creates an

208    entry in DataCite's metadata store. A second request is then immediately posted to the

---

[11] https://mds.datacite.org/static/apidoc

209    same API which registers the DOI and specifies its landing page (the page in the CEH

210    catalogue from which the administrator made the request). Next, a request is made to the

211    shortDOI service[12]  which creates a more practical, shorter DOI alias. Both the new DOI and

212    the shortDOI are then automatically added to the metadata record in the data catalogue,

213    along with information about how to cite the data resource (Fig 1). Once a DOI has been

214    registered for a data resource, subsequent updates or amendments to the metadata which

215    affect the DataCite metadata are automatically submitted to the DataCite API. This ensures

216    that the DataCite metadata is always representative of the GEMINI metadata held in the

217    CEH catalogue.

218

219    The researcher who deposited the data is emailed to inform them that a DOI has been given

220    to the data they created. The email contains details of the DOI, the shortDOI and

221    recommendations on how to use the DOI and cite the data. This notification is currently

222    carried out manually by a member of staff at the data centre. The EIDC also maintains an

223    inventory of all the datasets it holds that have a DOI. This DOI inventory is also manually

224    updated upon the registration of a new DOI. Whilst both these actions are currently carried

225    out manually, the EIDC hopes to automate them in future in order to reduce staff time spent

226    carrying out the processes and provide a more efficient service to depositors.

227

228    To date, just over 70% of the data resources held by EIDC have a DOI allocated to them.

229    Currently, researchers are asked upon deposit whether they would like a DOI for the data

230    they have created – they are not minted automatically for every data resource taken in. The

231    reason researchers don't always request a DOI is usually due to the data being 'legacy' data

232    i.e. data that was generated a long time ago (on the scale of decades) and has already been

233    discussed in the scientific literature, therefore researchers feel they have nothing to gain

---

[12] http://shortdoi.org/

234    from obtaining a DOI for them. When a DOI for a data resource is resolved using a web

235    browser, the user sees a landing page which is the discovery metadata record for that

236    resource. The landing page includes information on how to obtain the resource and how to

237    cite it in future publications (see above). DOIs can only be allocated to data resources that

238    have been formally deposited with the EIDC; this normally occurs towards the end of a

239    project or section of work. Data must have passed documented quality checks and be held

240    within the data centre itself. DOIs are allocated prior to the data being made publicly

241    available (although this usually happens immediately after). The EIDC supports NERC's

242    option of allowing researchers a two year embargo on the release of the data they created.

243    In the case of embargoed data resources, DOIs are registered when the data are deposited,

244    as this allows researchers to use the DOI in any publications they have planned. The DOI is

245    documented on the landing page for the data resource along with details of the embargo and

246    a date when the data are to be made available.

247

248    3.0 Uptake and use of DOIs for data from a data centre perspective

249    The motivation for requesting a DOI for data deposited with EIDC has varied over time. At

250    first, requests came in solely because it was now a service offered by the data centre and

251    this had been communicated to depositors by the IL staff. DOIs were initially requested even

252    though some researchers were not fully aware of what they could be used for. This is not

253    unsurprising, as it has been noted previously that there is a lack of clear recommendations

254    on how to cite data within scientific literature. The Data Citation Guidelines for Data

255    Providers and Archives [10] state that among Federation of Earth Science Information

256    Partner (ESIP) members, current recommendations for citing data range from casual

257    acknowledgement within the text of a paper to formal and specific citations within the

258    references section of the paper. Mayernik [5] also stated that even when data is widely

259    shared, users do not commonly cite datasets in formal ways. Rather than formally citing

260 datasets, data users typically acknowledge data use in the text of an article in the

261 acknowledgement section.

262

263 One of the first DOIs assigned by the EIDC was for data created by Beresford et al [13]

264 which was subsequently quoted in a journal paper [14]. However, the authors failed to

265 include the recommended DataCite citation in the reference list and merely added a

266 statement to the paper, "All data associated with this study are available from the CEH

267 Information Gateway (https://gateway.ceh.ac.uk/) and the data have been allocated a digital

268 object identifier (http://dx.doi.org/10.5285/1a91c7d1-ec44-4858-9af2-98d80f169bbd)"

269 This indicates they did not regard it as a reference in the same way as they would a journal

270 paper.

271

272 Other researchers requested a DOI as they were publishing in a data journal and it was a

273 mandatory requirement of submission. Data journals, especially in the field of the

274 environmental or natural sciences are a relatively new concept, however, they are increasing

275 in number. Journals such as Earth System Science Data (ESSD), Geoscience Data Journal,

276 Scientific Data and Data in Brief publish peer-reviewed data papers – papers that describe

277 datasets [3]. The majority of data journals require data to be stored in an approved

278 repository with a permanent identifier assigned enabling reviewers to access the data. At the

279 EIDC, one of the first datasets referred to in a data paper was from Haxton et al. [15]. This

280 dataset was deposited with the EIDC and given a DOI which was subsequently cited in an

281 ESSD paper by Prudhomme et al. [16]. Furthermore, the ESSD paper has since been cited

282 by at least five other journal publications (as recorded by CrossRef[13]) including one the

283 author co-authored [17]. It should be noted that each of these outputs (the data paper and

---

[13] http://www.crossref.org/

284  the dataset itself) is a publication in its own right - there is no requirement for the data

285  resource to have the same lead author as the data paper. They are separate entities with

286  their own individual reference and can be referred to as such. When re-using data or

287  tracking where data have been re-used it is important to use the citation for the dataset itself,

288  rather than the reference for the data paper. If simply referring to the work carried out by a

289  group of authors, citing the data paper would be appropriate. By publishing a data paper

290  based on a dataset, authors are adding value to the dataset for the future consumers of the

291  data [9] as the data they created has undergone a scientific peer review process. Datasets

292  published by the data centre have reached a certain level of quality as required to obtain a

293  DOI, but they are not peer reviewed.

294

295  As a case study, the above example has encouraged other researchers within the

296  organisation to engage with the data centre which has further increased the number of

297  datasets being offered for deposit. In the financial year 2012-2013, the EIDC had 35 deposit

298  requests i.e. researchers contacting the data centre wishing to deposit data. These figures

299  contrast with those of the financial year 2014-2015, where the EIDC had 83 deposit requests

300  (it should be noted that one deposit request may lead to the deposit of one dataset, or many

301  which is often the case). In 2015-2016, the EIDC had 61 deposit requests in the first 6

302  months of the financial year, therefore it is likely that the number of deposit requests this

303  year will exceed those of the previous year. The reason for this increase in engagement with

304  the EIDC could be due to case-studies such as the one above being advertised keenly

305  throughout the organisation (CEH), however, it is more likely that pressure from publishing

306  houses, as discussed below, has had a greater impact on these figures.

307

308  The final reason researchers are now offering their data to the data centre to publish (and

309  requesting a DOI) is that increasingly scientific journals are recommending, or even

310  mandating, that data referred to in an article must be archived in an appropriate public

311  archive [3]. The archive must provide public access and guarantee long-term preservation of

312  the data resource. Some journals also require that the data have been assigned a

313  permanent identifier (e.g. a DOI). The pressure from publishing houses (e.g. British

314  Ecological Society, Ecological Society of America, Nature and Science) is urging those

315  researchers in the environmental science community previously resistant to the idea of

316  publishing data to actively participate. Whilst this is encouraging it is often done in an

317  untimely manner. Despite the support and advice available, some researchers are still

318  unaware of the importance of data management and citation, or it fails to make the list of

319  their priorities for reasons discussed above [3].  Many researchers are currently offering data

320  resources to the data centre for publication only after a journal paper has been written and

321  submitted, and hence require the deposit process to take place hurriedly. This is often not

322  possible as the EIDC processes mandate that data coming into the data centre be

323  accompanied by sufficient supporting information which depositors have usually not

324  prepared in advance. The EIDC is bound by NERC to take in data of long-term value so that

325  it may be stored securely in perpetuity and have the potential to be re-used where suitable. It

326  is therefore not possible for the data centre to 'fast-track' data deposits with the aim of

327  meeting the requirement from depositors that they must have a DOI for data referred to in a

328  journal paper. Data accepted into the EIDC must be complete, be in a non-proprietary format

329  and have sufficient supporting information so that it may be understood and re-used by

330  others without the need to contact the creator. It is therefore critical that researchers engage

331  with data centre staff as early as practically possible in their projects, to develop data

332  management plans and ensure the correct documentation will be provided upon deposit of

333  the data. In cases where researchers have taken advantage of the support provided and

334  deposit of data has occurred in a timely manner, the process of obtaining a DOI and

335  publishing the data can occur rapidly as the workflow operated within the datacentre is

336  automated, where appropriate, and can be completed in a matter of seconds. If researchers

337  have not planned in advance and approach the data centre requesting a DOI as a matter of

338 urgency, the process can take somewhat longer. This is because time has to be spent

339 preparing the data and supporting information. Therefore, whilst the pressure from

340 publishing houses has prompted increased awareness of the requirement to publish data, it

341 may take some time before researchers realise they must engage with this process at an

342 early stage, before a project or grant is completed and prior to preparing articles for

343 submission.

344

345 Since the EIDC began issuing DOIs for data resources we have seen an increase in

346 researchers' awareness of the requirement to make data available, predominantly driven by

347 data journals and journal publications. The EIDC is receiving an increasing number of

348 enquiries about depositing data from scientists interested in submitting data papers and

349 research articles as they realise that this is a mechanism whereby they can gain academic

350 credit for a body of work which was previously unacknowledged. Our ability to identify and

351 cite data resources in a reliable manner is largely down to the system put in place by

352 DataCite and the use of DOIs (although it is possible that other permanent identifiers could

353 work in an equally successful way [8]) as it offers researchers an incentive for releasing the

354 data they have created. Without this incentive, we believe many data resources available

355 today through the EIDC would not have been deposited with the data centre and therefore

356 be inaccessible.

357

358 4.0 Future challenges for the data centre

359 The advent of a robust method for making data resources citeable has gone some way in

360 addressing the lack of published data available in the field of environmental sciences but

361 there are still areas where improvements could be made to further increase openness and

362 re-use of data. Many of the data resources archived by the EIDC are created from long-term

363 environmental monitoring programmes and therefore data are being regularly updated. The

challenge of making this type of dynamic dataset citeable is well documented, as data such as these do not fit the commonly used DOI system well [5, 11, 12]. In line with DataCite recommendations, once a dataset held by the EIDC has been given a DOI, it will not be changed, updated or corrected [18]. If any of these alterations are required, a new DOI is issued. This is so users can identify and retrieve the exact same data identified by a DOI irrespective of how long it has been since it was registered. The EIDC currently offers researchers two choices when depositing dynamic data, based on the approaches outlined by Ball and Duke [12]; either a new time-slice can be deposited into the data centre and a new DOI issued, or the whole dataset can be taken in including the previous data and any new data (a new snap-shot). In the latter case, the previous version is deprecated and a new DOI is issued for the whole resource. An example of this is data from the UK Butterfly Monitoring Scheme (UKBMS) deposited into the EIDC. The UKBMS deposits data annually on collated indices and species trends. The first deposit was made in 2011 and the data ran from 1976 to 2011 [19, 20]. In 2012, UKBMS submitted a new snap-shot of the data, this time running from 1976 to 2012 [21, 22]. The addition of the new data not only added extra data values but as a consequence also changed the values of the previous years' data. Once a new snap-shot is published and has a DOI, the old snap-shot is deprecated by labelling it an 'Historical archive' in the discovery metadata record. The catalogue is configured so that for records labelled as such, a banner automatically appears at the top of the record stating 'This dataset has been withdrawn' (Fig 2). In this way, the DOI still resolves to the correct landing page so remains a permanent identifier and the user can clearly see that this is not the most current version of the dataset (a link to a record for this collection of data resources is available from the deprecated dataset landing page so users can easily find the most up-to-date version, should they wish to).

However, some researchers are unhappy with the current system and indeed, from a data centre's perspective, snap-shots can become unwieldy for regularly updated time-series

391   data which are common in the Earth Sciences [10]. Instead, researchers would prefer one

392   identifier for the whole resource that never changed regardless of how many updates or

393   additions of data were made. Such a system would ensure citation metrics for the resource

394   were not diluted with new citations generated each time an update was made. However, this

395   is a service we are currently unable to offer using the system we have in place. The

396   Research Data Alliance[14] has a working group dedicated to exploring solutions to the

397   problem of citing dynamic datasets and a position paper by Andreas Rauber and Stephan

398   Pröll has been produced describing a conceptual model for scalable dynamic data citation

399   [23]. However, this paper addresses the problem from a data re-user's perspective so may

400   not solve the issues that researchers depositing to the EIDC have raised. Rauber and Pröll

401   propose using timestamped, versioned data that can be assembled into specific subsets by

402   using queries which subsequently have permanent identifiers assigned to them. This system

403   enables authors to cite only the query, rather than the whole dataset, ensuring users can

404   access exactly the same data referred to by the identifier for perpetuity [23].Whilst this

405   addresses the issue of ensuring users are able to precisely identify specific subsets of data

406   that may have changed over time it does not solve the issue of citation dilution raised by

407   researchers depositing to EIDC. Also, the DOI system, as currently implemented by

408   DataCite, does not support Template Handles, thus a parameterized DOI would not resolve

409   to a particular subset but to the whole dataset [11]. It is clear that attempts are being made

410   to address the issue of citing dynamic datasets but also that one size does not fit all [3, 5],

411   therefore systems may have to adapt in future to accommodate researchers requirements.

412

413   Another issue which, if addressed, could further promote data publication in the

414   environmental sciences is that of citation metrics. For the production and publication of data

415   to be recognised as valuable scholarship it requires accepted metrics of significance [7].

---

[14] https://rd-alliance.org/

416    Researchers are more likely to publish the data they create if they can measure its impact,

417    track its use and receive credit for creating it [3, 7]. A researcher's academic success is

418    frequently measured by the journal publications they produce, specifically in the impact

419    factor of the journals in which they publish and the number of times articles are subsequently

420    cited. If mechanisms were put in place to provide similar information for datasets,

421    researchers would be able to measure the impact of the data they produce which could input

422    into the professional reward process [5]. Tracking data use is difficult as datasets are

423    inconsistently cited by data users [5]. However, respondents to a survey carried out by Kratz

424    and Strasser [7] found that citation and download counts were more useful than search rank

425    or altmetrics. Therefore, a method for measuring data impact based on data citation counts,

426    though difficult to implement would be desirable to researchers. Data papers can go some

427    way to providing this type of information. For example, the data journal ESSD provides

428    metrics on views and citations of the data papers they publish (Fig. 3). Crucially, however,

429    this is not tracking the citation of the data itself which has its own DOI and mechanism for

430    citation. Thomson Reuters[15] Web of Knowledge now provide a service called the Data

431    Citation Index (DCI) which provides access to data, links data to the articles it supports and

432    tracks citation of datasets. Unfortunately, the DCI is currently not open and free to use (a

433    subscription is required) and repositories have to agree to have information about the data

434    they hold harvested by Thomson Reuters. The EIDC is working with Thomson Reuters to

435    ensure that the data it holds can be included in the DCI and this has recently been achieved

436    by allowing Thomson Reuters to harvest metadata held by DataCite about data held by the

437    EIDC. This is an important first step, although, as CEH is not a subscriber to the DCI, the

438    data centre is unable to obtain information on the citation counts for data it holds.

439

---

[15] http://thomsonreuters.com/en.html

440    In contrast, ResearchGate is a free service, which enables researchers to share their

441    publications and access citation metrics. ResearchGate allows registered users to add

442    articles, book chapters, conference papers, datasets and unpublished work to their home

443    page - once added, metrics are collected on the publication. This would seem like a suitable

444    solution to the problem of collecting citation metrics for the data held by the EIDC, at least

445    from an individual researcher's point of view, as they should theoretically be able to include

446    information about datasets they have deposited with the datacentre (e.g. title, DOI) and

447    obtain information of citation metrics over time. However, when registering a dataset with

448    ResearchGate, users are required to attach the data as a file and are therefore uploading a

449    copy of the data to the ResearchGate site. This is not something the EIDC can recommend

450    for a number of reasons. First, additional copies of data would be unnecessarily generated

451    and stored. Second, ResearchGate mandate that any data uploaded to its site is free from

452    any Intellectual Property Rights which in the majority of cases is not true for data generated

453    though public or private funding. Third, uploading data to ResearchGate is often impractical

454    as the volumes of data in question are often very large (500GB or more)

455

456    It is clear that whilst some solutions are available, further work is still needed to implement

457    an openly accessible tool to capture and present metrics for datasets. Until researchers can

458    quantify the impact data resources they have created have on the academic community as a

459    whole they may not receive the full scholarly credit they deserve. In the meantime, the EIDC

460    plans to include information on download counts for each dataset on its landing page. Whilst

461    not ideal, it provides researchers with a highly regarded 'second-choice metric' [7] and can

462    be used as an interim measure until a more informative system is put in place.

463

464    5.0 Conclusions

465 Whilst there is still a long way to go before data resources are viewed as valued outputs

466 from a researcher's work in the same way journal publications have always been, data

467 centres, such as the EIDC, are facilitating a cultural shift in practices with regard to data

468 publications. By providing a robust workflow enabling the identification of datasets and

469 providing a means for data to be cited, data centres are providing the building blocks on

470 which more wholesale changes in attitude and behaviour can occur. Working in conjunction

471 with publishing houses, data centres are beginning to convince researchers that publishing

472 the data they have generated can be beneficial to their research careers. Data centres can

473 further improve on the volume of data published in the environmental sciences by enabling

474 the citation of dynamic datasets, ensuring long-term environmental monitoring experiments

475 can be cited as a single entity, rather than having to generate a new DOI and citation after

476 each new addition of data. In addition, the generation and publication of citation metrics that

477 provide an indication of the impact a dataset has had on the academic community could

478 also, encourage more researchers to publish the data they have created. Much has been

479 accomplished in the last few years but there are still many issues left to address. It will take

480 time for a cultural shift to occur, but by putting flexible robust systems in place and by

481 seeking to illustrate to researchers the benefits of publishing the data they produce, in time

482 data resources and those that generate them will receive the credit and standing they

483 deserve.

484

485 Acknowledgements

489

490 References

491

492    1. Lawrence B, Jones C, Matthews B, Pepler S, Callaghan S (2011) Citation and peer review

493    of data: Moving towards formal data publication. Int J Digital Curation. 6(2): 4-37

494

495    2. Klump J, Bertelmann R, Brase J, Diepenbroek M, Grobe H, Hock H, Lautenschalger M,

496    Schindler U, Sens I, Wachter J (2006) Data publication in the Open Access Initiative. Data

497    Science Journal 5: 79-83.

498

499    3. CODATA-ICSTI. (2013) ed. Y Socha. Out of Cite, Out of Mind: The Current State of

500    Practice, Policy, and Technology for the Citation of Data. Data Science Journal, 12(0):

501    CIDCR1-CIDCR7. DOI: http://doi.org/10.2481/dsj.OSOM13-043

502

503    4. Arzberger P, Schroeder P, Beaulieu A, Bowker G, Casey K, Laaksonen L, Moorman D,

504    Uhlir P, Wouters P (2004) Promoting access to public research data for scientific, economic,

505    and social development. Data Science Journal 3: 135-152

506

507    5. Mayernik, M, (2013) Bridging data lifecycles: Tracking data use via data citations

508    workshop report. NCAR Technical Note NCAR/TN-494+PROC,

509    http://dx.doi.org/10.5065/D6PZ56TX.

510

511    6. Assante M, Candela L, Castelli D,  Manghi P, Pagano P (2015) Science 2.0 repositories:

512    Time for a change in scholarly communication. D-Lib Magazine 21:(1/2)

513    doi:10.1045/january2015-assante. http://dx.doi.org/10.1045/january2015-assante

514

515     7. Kratz JE, Strasser C (2015) Making Data Count. Scientific Data 2:150039. Doi:

516     10.1038/sdata.2015.39

517

518     8. Duerr R, Downs R, Tilmes C, Barkstrom B, Lenhardt W, Glassy J, Bermudez L,

519     Slaughter P (2011) On the utility of identification schemes for digital earth science data: An

520     assessment and recommendations. Earth Science Informatics 4:139-60.

521     http://dx.doi.org/10.1007/s12145-011-0083-6

522

523     9. Callaghan S, Donegan S, Pepler S,*et al.* (2012) Making data a first class scientific output:

524     Data citation and publication by NERC's environmental data centres. Int J Digital Curation.

525     7(1): 107–113.

526

527     10. ESIP (Federation of Earth Science Information Partners) (2012) Data Citation Guidelines

528     for Data Providers and Archives. Ed. Parsons MA, Barkstrom B, Downs RR, Duerr R,

529     Tilmes C and ESIP Preservation and Stewardship Committee. ESIP Commons.

530     http://dx.doi.org/10.7269/P34F1NNJ

531

532     11. Klump J, Huber R, Diepenbroek M (2015) DOI for geoscience data-how early practices

533     shape present perceptions. Earth Science Informatics 1-14.

534     http://dx.doi.org/10.1007/s12145-015-0231-5

535

536     12. Ball A, Duke M (2015) How to Cite Datasets and Link to Publications. DCC How-to

537     Guides. Edinburgh: Digital Curation Centre. http://www.dcc.ac.uk/resources/how-guides

538

539  13. Beresford NA, Barnett CL, Howard BJ, Howard DC, Tyler AN, Bradley S, Copplestone D

540  (2011). Observations of Fukushima fallout in Great Britain. NERC Environmental Information

541  Data Centre. doi: 10.5285/1a91c7d1-ec44-4858-9af2-98d80f169bbd

542

543  14. Beresford NA, Barnett CL, Howard BJ, Howard DC, Wells C, Tyler AN, Bradley S,

544  Copplestone D (2012) Observations of Fukushima fallout in Great Britain. J Environ

545  Radioact. 114:48-53. doi: 10.1016/j.jenvrad.2011.12.008.

546

547  15. Haxton T, Crooks S, Jackson CR, Barkwith AKAP, Kelvin J, Williamson J, Mackay JD,

548  Wang L, Davies, H, Young A, Prudhomme C (2012). Future flows hydrology data. NERC

549  Environmental Information Data Centre. doi:10.5285/f3723162-4fed-4d9d-92c6-

550  dd17412fa37b

551

552  16. Prudhomme C, Haxton T, Crooks S, Jackson C, Barkwith A, Williamson J, Kelvin J,

553  Mackay J. Wang L, Young A, Watts G (2013) Future Flows Hydrology: and ensemble of a

554  daily river flow and monthly groundwater levels for use for climate change impact

555  assessment across Great Britain. Eath Syst. Sci. Data 5:101-107. doi:10.5194/essd-5-101-

556  2013

557

558  17. Royan A, Prudhomme C, Hannah DM, Reynolds SJ, Noble DG, Sadler JP (2015)

559  Climate-induced changes in river flow regimes will alter future bird distributions. Ecosphere

560  6(4): 50. doi: 10.1890/ES14-00245.1

561

562    18. British Library, DataCite (2013) Working with the British Library and DataCite – A guide

563    for Higher Education Institutions in the UK. British Library

564    http://www.bl.uk/aboutus/stratpolprog/digi/datasets/WorkingWithDataCite_2013.pdf

565

566    19. Botham, M; Roy, D; Brereton, T; Middlebrook, I; Randle, Z (2012). United Kingdom

567    Butterfly Monitoring Scheme: collated indices 2011. NERC Environmental Information Data

568    Centre.doi:10.5285/ff55462e-38a4-4f30-b562-f82ff263d9c3

569

570    20. Botham, M; Roy, D; Brereton, T; Middlebrook, I; Randle, Z (2013). United Kingdom

571    Butterfly Monitoring Scheme: species trends 2011. NERC Environmental Information Data

572    Centre.doi:10.5285/cad2af6c-0c97-414c-8d5f-992741b283cf

573    20

574

575    21. Botham, M.; Roy, D.; Brereton, T.; Middlebrook, I.; Randle, Z. (2013). United Kingdom

576    Butterfly Monitoring Scheme: collated indices 2012. NERC Environmental Information Data

577    Centre.doi:10.5285/7949cc99-76c4-4a3e-8c33-41a35b8b7777

578

579    22. Botham, M.; Roy, D.; Brereton, T.; Middlebrook, I.; Randle, Z. (2013). United Kingdom

580    Butterfly Monitoring Scheme: species trends 2012. NERC Environmental Information Data

581    Centre.doi:10.5285/5afbbd36-2c63-4aa1-8177-695bed98d7a9

582

583    23. Rauber A, Pröll S (2015) Scalable dynamic data citation – RDA-WG-DC Position paper

584    https://rd-alliance.org/groups/data-citation-wg/wiki/scalable-dynamic-data-citation-rda-wg-dc-

585    position-paper.html Accessed 25 June 2015

586

587

588

589

**Figure captions**

**Fig 1** Example of a record in the CEH catalogue showing recommended citation and display of DOI.

**Fig 2** Example of a deprecated metadata record in the CEH catalogue.

**Fig 3** Metrics provided by the data journal Earth System Science Data including views and citations.

597

598

599 **Table 1.** Mapping between GEMINI metadata and DataCite metadata

| GEMINI metadata element | DataCite metadata element |
|---|---|
| Based on /MD_Metadata/fileIdentifier | /resource/identifier |
| /MD_Metadata/identificationInfo/MD_DataIdentification/pointOfContact/CI_ResponsibleParty[role/CI_RoleCode/@codeListValue='author']/individualName | /resource/creators/creator |
| /MD_Metadata/identificationInfo/MD_DataIdentification/citation/CI_Citation/title | /resource/titles/title |
| /MD_Metadata/identificationInfo/MD_DataIdentification/pointOfContact/CI_ResponsibleParty[role/CI_RoleCode/@codeListValue=publisher]/individualName | /resource/publisher |
| /MD_Metadata/identificationInfo/MD_DataIdentification/citation/CI_Citation/date/CI_Date[dateType/CI_DateTypeCode/@codeListValue='publication']/date | /resource/publicationYear |
| /MD_Metadata/identificationInfo/MD_DataIdentification/descriptiveKeywords/MD_Keywords/keyword | /resource/subjects/subject |
| - | /resource/dates/date[@dateType='Submitted'] |
| /MD_Metadata/identificationInfo/MD_DataIdentification/language/LanguageCode | /resource/language |
| /MD_Metadata/MD_ScopeCode/@codeListValue | /resource/resourceType/@resourceTypeGeneral |
| /MD_Metadata/identificationInfo/MD_DataIdentification/citation/CI_Citation/identifier/RS_Identifier | /resource/alternateIdentifiers/alternateIdentifier |
| /MD_Metadata/distributionInfo/MD_Distribution/distributionFormat/MD_Format/name | /resource/formats/format |
| /MD_Metadata/identificationInfo/MD_DataIdentification/resourceConstraints/MD_LegalConstraints/otherConstraints | /resource/rightsList/rights |
| /MD_Metadata/identificationInfo/MD_DataIdentification/abstract | /resource/descriptions/description[@descriptionType='Abstract'] |
| /MD_Metadata/identificationInfo/MD_DataIdentification/extent/EX_Extent/geographicElement/EX_GeographicBoundingBox | /resource/geoLocations/geoLocation/geoLocationBox |

600

Figure

Dataset
# Biomass of Trifolium repens versus Lolium perenne after ozone exposure in solardomes

Hayes, F.; Mills, G.; Ashmore, M. (2014)

doi:10.5285/5de90f7a-dec9-4bd5-af53-d0873a89d25d

The data are biomass measurements from an ozone exposure experiment, during which Trifolium repens and Lolium perenne were exposed as both monocultures and two-species mixtures to an episodic rural ozone regime in large, well-watered containers within solardomes for 12 weeks. Treatments were elevated ozone (AOT40 (Accumulated Ozone Threshold exposure of 40 parts per billion) of 12.86 ppm h) or control conditions (AOT40 of 0.02 ppm h). Measurements were dry weight, with a cutting height of 7cm above soil level. The distribution of plant material within the canopy was determined by separating material growing in the upper canopy (>14cm) from the canopy edge and the inner canopy for both species. The experiments were carried out in the CEH Bangor Air Pollution Facility. Work was funded by the Centre for Ecology and Hydrology Integrating Fund Initiative. The observed decreases in photosynthetic efficiency and capacity in elevated ozone indicate that the ability of such ubiquitous vegetation to act as a sink for atmospheric carbon may be reduced in future climates.

Publication date: 2014-12-18 ( created 2010-01-01 )

## Where/When

**Study area**

**Temporal extent**
2007-04-30   to   2007-10-31

## Online Resources

Link to paper on NERC Open Research Archive (NORA)
Hayes, Felicity; Mills, Gina; Ashmore, Mike. 2010 How much does the presence of a competitor modify the within-canopy distribution of ozone-induced senescence and visible injury? Water, Air and Soil Pollution, 210. 265-276. 10.1007/s11270-009-0248-9

Link to paper on NERC Open Research Archive (NORA)
Hayes,F.; Mills, G.; Ashmore, M.. 2009 Effects of ozone on inter- and intra-species competition and photosynthesis in mesocosms of Lolium perenne and Trifolium repens. Environmental Pollution, 157 (1). 208-214. 10.1016/j.envpol.2008.07.002

Supporting information
Supporting information available to assist in re-use of this dataset.

Online ordering
Order a copy of this database

### Get the data

⬇ Online ordering

📄 Supporting documentation

Format of the data: Comma Separated Values

**If you reuse this data, you must cite**
Hayes, F.; Mills, G.; Ashmore, M. (2014). Biomass of Trifolium repens versus Lolium perenne after ozone exposure in solardomes. NERC Environmental Information Data Centre. doi:10.5285/5de90f7a-dec9-4bd5-af53-d0873a89d25d

BibTeX   RIS

This resource is made available under the terms of the Open Government Licence

### This dataset is part of the series

Trifolium and Lolium competition in ozone

Figure

Figure

Earth System Science Data
The Data Publishing Journal

| Imprint | Contact |

Volume 1, Issue 1

Article  Metrics  Related Articles

Search ESSD

Final revised paper

13 Mar 2013

## Future Flows Hydrology: an ensemble of daily river flow and monthly groundwater levels for use for climate change impact assessment across Great Britain

C. Prudhomme, T. Haxton, S. Crooks, C. Jackson, A. Barkwith, I. Williamson, J. Wallis, J. Mackay, L. Wang, A. Young, and G. Watts

### Viewed

Total article views: 1,488 (including HTML, PDF, and XML)

| HTML | PDF | XML | Total | BibTeX | EndNote |
|------|-----|-----|-------|--------|---------|
| 741 | 690 | 57 | 1,488 | 40 | 11 |

Views and downloads (calculated since 13 Mar 2013, article published on 13 Mar 2013)



HTML views   PDF downloads   XML downloads

Cumulative views and downloads (calculated since 13 Mar 2013, article published on 13 Mar 2013)



HTML views   PDF downloads   XML downloads

### Cited

crossref  Google Search
5

### Saved

citeulike   Mendeley
1   15

### Discussed

Latest update: 11 Nov 2013 22:39