# Toward Reliable Biodiversity Data References

Michael Elliott[1], Jorrit H. Poelen[2,*], José A.B. Fortes[1]

**1 Advanced Computing and Information Systems Laboratory (ACIS), University of Florida, Gainesville, Florida, USA**
**2 400 Perkins St Apt 104, Oakland, California, USA**

**\* jhpoelen@jhpoelen.nl**

## Abstract

Scientific discovery increasingly relies on digital datasets to capture measurements and outcomes. However, no systematic approach has been adopted to reliably reference and provide access to our digital datasets. Our existing data infrastructures have grown accustomed to using location-based identifiers such as URLs in an attempt to retain our digital knowledge. We hypothesize that URLs are not sufficient to ensure long-term data access, then propose a method for evaluating long-term URL reliability.

After taking periodic inventories from March through October 2019 of the data served by major biodiversity aggregators, including GBIF, iDigBio, DataONE, and BHL, we found that, for each network, 5%-44% of registered URLs were intermittently or consistently unresponsive, 0%-64% produced unstable content, and 13%-76% became either unresponsive or unstable over the period of

observation. We propose to use content-based identifiers to reliably track and reference datasets while enabling decentralized archiving schemes. We propose a method for properly tracking and archiving datasets that can be used to guarantee fixed content and encourage long-term accessibility by leveraging content- rather than location-based identifiers.

# Introduction

Over the course of hundreds of years, naturalists and biologists have systematically collected physical evidence from an ever-changing natural world. Through well-established protocols and institutional support, many of these natural history collections have withstood the ravages of time [Hortal et al., 2015, Davis and Schmidt, 1996]. Records that describe these carefully collected specimens are now made available digitally through online search indices, registries, and data archives [Page et al., 2015]. The increased availability of digital natural history records helps work toward Charles Elton's realization that ecosystems can only be fully understood when we "provide conceptions which can link up into some complete scheme the colossal store of facts about natural history which has accumulated up to date in this rather haphazard manner" [Elton, 1927]. So far, various initiatives have succeeded to provide comprehensive aggregate views from previously scattered natural history record siloes [Rinaldo and Norton, 2009, Michener et al., 2011, Edwards, 2000, Matsunaga et al., 2013, Facility, 2019]. However, we show that these aggregate views are

subject to change as their underlying digital source data changes or becomes inaccessible. Although efforts have been made to keep track of changes in digital networked resources, such as the use of version numbers and last modified dates [Wieczorek et al., 2012, Robertson et al., 2014] and periodic archival [Costello et al., 2013], we are not aware of the adoption of any systematic approach to preserve the accessibility as well as longevity of our digital natural history record and derived datasets. We have collected evidence that, despite hundreds of years of experience in preserving our physical natural history records, we are currently faced with a growing body of digital data that changes daily and can disappear with the push of a button. Our scholarly record is stitched together by an intricate web of associations between scientific publications. These associations are made explicit using citations. These citations point to related scientific works and are assumed to provide enough identifying information to allow the reader to retrieve the unaltered referenced work regardless of the time at which the reader chooses to do so [Garfield et al., 1964]. In the pre-internet era, the lookup of these references required access to one of the many academic libraries in the world. With the rise of internet accessible scientific publications, authors and readers access these references using a networked device by downloading content from publication websites. This means that researchers are increasingly citing online works to support their claims. Because the citation format of online works documents only when (e.g., 2019-10-01) and where (e.g., https://doi.org/10.123/456) the referenced work was accessed by the author [GBIF.org, 2019, iDigBio.org, 2016, DataONE, 2012], the future reader expects the web accessed resource to remain accessible and unaltered via this single web location. Future

readers may attempt to find a version of the works referenced by searching 67
online data networks for the matching author and title, but there is no 68
guarantee that information found this way will be exactly the same as 69
what was originally referenced. Any reference that does not allow future 70
readers to find the referenced work fails to satisfy the FAIR principle of 71
findability: "F1. (meta)data are assigned a globally unique and eternally 72
persistent identifier." [Wilkinson et al., 2016]. Our study is not alone in 73
providing evidence that suggests that networked, location-based access to 74
digital objects is an unreliable mechanism for providing continued access to 75
the unaltered original work [Vision, 2010, Klein et al., 2014]. Unless we 76
change the way we preserve and cite our digital scholarly works, our 77
physical records stored in libraries and museums around the world are 78
likely to outlast our digital ones. 79

## Problem Characterization 80

We show that the current practice of using Uniform Resource Locators 81
(URLs) [Berners-Lee et al., 1994] to reference online biodiversity datasets 82
provides no guarantee of long-term data accessibility. Readers who 83
encounter references that use URLs as dataset identifiers cannot be certain 84
that the referenced data will continue to be accessible and in its exact 85
original form. This uncertainty might be cause for alarm for researchers 86
because, over time, the integrity of the scholarly record itself is damaged 87
when existing references become reliable due to the loss of access to the 88
data they reference. When data access is lost, it is possible that 89
documented research results may become impossible to reproduce and the 90
justification for any conclusions or hypotheses that relied on lost results 91

4

may be undermined. If the use of error-prone referencing techniques is not <sub>92</sub> addressed, we expect that any resulting gaps in the biodiversity data <sub>93</sub> record will only become more severe. <sub>94</sub>

The current practice of relying on URLs to locate and identify <sub>95</sub> referenced data is hazardous due to their demonstrated risk of link rot and <sub>96</sub> content drift [Klein et al., 2014]. Link rot occurs when a URL, or link, that <sub>97</sub> had previously responded to queries can no longer be reached. This can <sub>98</sub> happen, for example, due to temporary outages, URL retirement, or URL <sub>99</sub> migration. A link exhibits content drift when a query to the link provides <sub>100</sub> content that is different from the content it provided in the past. The <sub>101</sub> extent of content drift can vary; content may have received only minor <sub>102</sub> edits with no changes in semantics, or it may reference a different entity <sub>103</sub> altogether. When a single URL is used to locate data that may change <sub>104</sub> over time, access to any particular version of the data is likely to be <sub>105</sub> short-lived. We show that, in the event of link rot or content drift, any <sub>106</sub> existing references that relied affected URL may become unreliable. <sub>107</sub>

In one study on the Genetics journal, it was reported that 40% of links <sub>108</sub> (URLs) to supplemental materials became unavailable due to link rot <sub>109</sub> within one year of publication [Vision, 2010]. Another study [Klein et al., <sub>110</sub> 2014] confirmed that as many as one in five articles in journal of Science, <sub>111</sub> Technology, and Medicine provide references that exhibit either link rot <sub>112</sub> and content drift and refer to the existence of either as "reference rot". <sub>113</sub> Since existing biodiversity references largely rely on URLs to locate <sub>114</sub> datasets, it is reasonable to expect that biodiversity data networks are also <sub>115</sub> at risk of providing unreliable dataset references as a result of reference rot. <sub>116</sub> The information systems used by major biodiversity data networks, such as <sub>117</sub>

DataONE, GBIF, and iDigBio, rely on data curators, such as institutional <sub>118</sub> repositories, to maintain active dataset URLs, and aggregate the data <sub>119</sub> found at those URLs for distribution in response to user queries. If a data <sub>120</sub> curator modifies, relocates, or stops serving a particular dataset, it may <sub>121</sub> become impossible to retrieve the original dataset and the integrity of the <sub>122</sub> data network will suffer as a result. <sub>123</sub>

In this paper, we propose a methodology for measuring the existence of <sub>124</sub> link rot and content drift in online data networks, then provide <sub>125</sub> experimental results that confirm the existence of both link rot and <sub>126</sub> content drift across all of the biodiversity data networks we considered, <sub>127</sub> including BHL, DataONE, iDigBio, and GBIF. Finally, we propose a <sub>128</sub> method for referencing and serving biodiversity data in a way that works <sub>129</sub> toward satisfying the Findable, Accessible, Interoperable, and Reusable <sub>130</sub> (FAIR) principles [Wilkinson et al., 2016]. <sub>131</sub>

## Methodology <sub>132</sub>

Although it has been demonstrated that reference rot does occur when <sub>133</sub> URLs are used for referencing scientific works [Vision, 2010, Klein et al., <sub>134</sub> 2014], we are not aware of any prior studies that provide quantitative <sub>135</sub> evidence that reference rot occurs specifically in biodiversity data networks. <sub>136</sub> We set out to quantify the extent of reference rot in biodiversity data <sub>137</sub> networks. Because reference rot occurs in the scope of individual data <sub>138</sub> references, and references to digital datasets rely on URLs to locate the <sub>139</sub> data, we begin by introducing terminology for characterizing the reliability <sub>140</sub> of a URL according to how often it exhibits link rot and content drift. <sub>141</sub>

6

## URL Reliability

We assume that the URLs used to reference biodiversity datasets are
expected to resolve to an Internet Protocol (IP) address in the Domain
Name System. If a web server exists at the resolved IP address, a query to
that address over the Hypertext Transfer Protocol (HTTP) will return a
response code and, in some cases, associated content [Berners-Lee et al.,
2005]. We classify the reliability of a URL according to the content, or lack
of it, that it provides over successive queries. If a query to a URL is
unsuccessful, we say that link rot has occurred. However, if a successful
response is received but the retrieved content is different from the content
retrieved by previous query, we say that content drift has occurred.
Monitoring URLs in this way allows us not only to determine whether link
rot and content drift occur, but also to capture their long-term behaviors.
For example, one URL that has exhibited link rot might have failed to
respond only once, whereas another might have become repeatedly
unresponsive. Likewise, one URL might exhibit content drift less frequently
than another whose contents change rapidly. Furthermore, various
combinations of link rot and content drift behavior may indicate that one
URL is more reliable than another, even though both exhibit reference rot.

We label URLs with sets of reliability indicators according to their link
rot and content drift behaviors. The defined reliability indicators are
differentiated by the degree of link rot and content drift observed over a
series of queries to the URL at different points in time. We characterize
the responsiveness of a URL according to how often it exhibits link rot:

- Unresponsive: the link has failed to respond to one or more queries

- Responsive: the link has responded to all recorded queries  167

We characterize the stability of a URL according to how often it  168
produces different content from one query to the next:  169

- Unstable: the content that the link points to sometimes changes  170

- Stable: the content that the link points to never changes  171

We characterize the overall reliability of a URL according to both of its  172
responsiveness and stability:  173

- Unreliable: the link does not always provide the expected content; it  174
  is either unresponsive, unstable, or both  175

- Reliable: the link always provides the expected content; it is both  176
  responsive and stable  177

Before we can determine the reliability of any given URL, we must first  178
monitor its behavior over time by documenting how it responds to periodic  179
queries. For the context of biodiversity, we consider the case when the  180
content that a URL produces is a dataset.  181

## The Data Collection Process  182

We suggest that digital dataset collection practices have some analogies to  183
well-established physical specimen collection procedures (Fig. 1) [Poelen,  184
2019g]. If datasets are considered analogous to specimens, then the URLs  185
that locate datasets are analogous to the physical locations of specimens in  186
the natural world; they are where digital datasets were originally found,  187
but not where they should be preserved. Once found, physical specimens  188

are collected by hand; similarly, digital datasets are downloaded by 189

querying their URLs. Once a specimen is collected and deposited to a safe, 190

well-known repository, a record is kept that documents what the specimen 191

is in addition to when, where, and by whom it was collected. 192

(insert figure 1 / see appendix) 193

The same can be done for downloaded datasets. When a dataset is 194

downloaded, a record can be kept that details the URL that was queried, 195

the time of query, and who (e.g. a human or software agent) issued the 196

query that initiated the download event; we refer to this record as the 197

dataset's provenance record. Additionally, the dataset itself should be 198

stored in a safe, well-known dataset archive. The final step in the 199

collection process is to link the actual preserved specimen to its 200

corresponding record (the "specimen history" in Fig. 1) via an assigned 201

unique identifier. For digital datasets, we use cryptographic hashes of the 202

data as unique content-based identifiers. 203

## Data Collection Over Time 204

By establishing a dedicated data observatory that follows the collection 205

process we have described, we can build a history for each observed URL 206

to capture its long-term reliability. Such an observatory should periodically 207

query the URLs listed in data network's URL registry, producing for each 208

URL two complementary parts: 1) an archived copy of the response to the 209

corresponding query, whether it was a dataset, an error code, or no reply 210

at all, and 2) a record of its provenance, including the URL itself, the 211

current date, and a content-based identifier of any dataset received. The 212

use of a content-based data identifier is crucial; it allows us to reliably link 213

9

each acquired dataset to its provenance record without the need for an    214
intermediate index. Successive provenance records can be aggregated to    215
construct comprehensive histories for both datasets (when and where they    216
were found) and URLs (which datasets they produced over a series of    217
queries over time).    218

The constructed URL histories can be analyzed to determine whether a    219
link was ever broken, when it was broken, and whether it became    220
responsive again. The logs also identify the content (or lack of it) that a    221
URL produced each time it was queried. Any change in the content    222
identifier from query to the next indicates a change in the content of the    223
dataset. These link breakages and content changes correlate to link rot and    224
content drift, respectively, and allow us to determine the responsiveness,    225
stability, and reliability of each URL over time.    226

## Data Network Reliability    227

Now that we have outlined a method for observing and documenting the    228
behavior of URLs over an extended period of time, we can apply our    229
method to observe all of URLs registered by biodiversity data networks.    230
We also extend the idea of URL reliability to entire data networks and    231
propose that the overall reliability of a data network can be evaluated by    232
monitoring the long-term reliability of each individual URL in the network    233
exposes. Whereas we rigidly label individual URLs with binary indicators    234
of responsiveness, stability, and reliability, we grade data networks    235
according to the percentage of registered URLs that are assigned each of    236
the reliability indicators. For example, if a data network contains three    237
distinct URLs and we find that only two out of the three are reliable, then    238

10

we say the data network is 67% reliable. <sub>239</sub>

## Experiment

The Preston biodiversity dataset tracker [Poelen et al., 2018] implements <sub>241</sub> mechanisms for monitoring data networks as we have described. It allows <sub>242</sub> users to deploy a data network observatory which systematically observes <sub>243</sub> the entire set of URLs registered by the network, queries each URL for <sub>244</sub> data, then documents data collection and archives the results. All crawl <sub>245</sub> activities, the queries they issue, and the results they produce are <sub>246</sub> meticulously recorded in a string of provenance logs. <sub>247</sub>

We deployed several Preston observatories which periodically queried <sub>248</sub> the registered dataset URLs listed by Biodiversity Heritage Library (BHL), <sub>249</sub> Data Observation Network for Earth (DataONE), Global Biodiversity <sub>250</sub> Information Facility (GBIF), and Integrated Digitized Bio Collections <sub>251</sub> (iDigBio). Each of these networks provides online registries of URLs that <sub>252</sub> locate the data in the network. The registered URLs for DataONE, GBIF, <sub>253</sub> and iDigBio were queried monthly from March 2019 through October 2019. <sub>254</sub> BHL was queried monthly from May 2019 through October 2019. The logs <sub>255</sub> taken by each of these observatories describe the URL queries and their <sub>256</sub> results, which were processed to produce the results that follow. A sixth <sub>257</sub> observatory was constructed by aggregating the queries of the five data <sub>258</sub> network observatories. <sub>259</sub>

(insert figure 2, see appendix) <sub>260</sub>

11

# Results

Breakdowns of the overall reliabilities of the data networks are provided in ²⁶²
Table 1. Results are listed as percentages and total counts of URLs in the ²⁶³
data network that were assigned each reliability indicator. When analyzing ²⁶⁴
the recorded results of queries to URLs in each data network over a period ²⁶⁵
of seven months, we found that, for each individual network, 5%-44% of ²⁶⁶
registered URLs were intermittently or consistently unresponsive, 0%-64% ²⁶⁷
produced unstable content, and 13%-76% became either unresponsive or ²⁶⁸
unstable over the period of observation. ²⁶⁹

Overall, 30% of URLs observed across the five networks became ²⁷⁰
unreliable at some point over the period of March 2019 through October ²⁷¹
2019. Of those unreliable URLs, 48% were unstable, 22% became ²⁷²
consistently unresponsive, and 70% were at best only intermittently ²⁷³
responsive. For 5% of successful queries, the URL failed to respond to the ²⁷⁴
next query. For 4% of successful queries, the URL provided different ²⁷⁵
content the next time it responded when queried. ²⁷⁶

The changes in reliability over time for each network are visualized in ²⁷⁷
Fig. 2. Note that because we have defined reliable URLs to be those ²⁷⁸
considered both responsive and stable, they always represent the smallest ²⁷⁹
fraction of URLs in Fig. 1, Fig. 2, and Fig. 3 visualizes the cumulative ²⁸⁰
growth of biodiversity data networks during their periods of observation. ²⁸¹
This growth is illustrated with two metrics: the total number of unique ²⁸²
URLs ever registered by each network and the total number of unique ²⁸³
contents that had been downloaded from the network at each sampled ²⁸⁴
point in time. ²⁸⁵

(insert figure 3, see appendix) ²⁸⁶

The behaviors of the distributions over time of responsive, stable, and reliable URLs vary notably between data networks.

(insert table 1, see appendix) Some reasons for these differences can be inferred when cross-examining the table and figures. For example, although BHL scored relatively low in responsiveness due to frequent link rot, the content that it does provide is more stable than all other networks because content drift within BHL is relatively rare. Conversely, although iDigBio is relatively responsive, it has low stability because the network's near-constant content growth far outpaces its URL growth. GBIF's behavior was characterized by large sporadic swings; a mass URL migration of over 14,000 Plazi-hosted datasets occurred in May, introducing thousands of new URLs over a short period of time, while over 31,000 URLs (60% of URLs that responded to queries that month) suddenly changed contents in October. Even the most reliable network, DataONE, shows a clear downward trend in all three categories, with 13% of URLs becoming unreliable over a period of just seven months. Additionally, DataONE's growth curves indicate that there are far fewer unique contents than unique URLs; this evokes two possibilities: either much of DataONE's URL population is unresponsive, or DataONE lists multiple URLs for many of its datasets. Because DataONE has been shown to be highly responsive, it could be the case that many distinct URLs refer to the same datasets. It's also worth noting that the June and September spikes in BHL's unresponsiveness were largely due to URLs that failed to respond in those particular months but actually did respond to future queries.

## Sources of Potential Numerical Error

We expect that the URL reliability counts generated for the figures and tables are lower than their actual values. When we qualified URLs as being reliable, responsive, and stable, we could not be certain that links did not briefly become unresponsive or change content during the month-long periods between queries. It is therefore likely that some cases of link rot and content drift were not reflected in the results. Additionally, we only query URLs that the data networks list in their dataset registries; this means that, after URL was removed from a network's registry, we could not detect subsequent instances of reference rot. Therefore, our results represent a very optimistic upper bound on URL and network reliabilities.

The results for DataONE and GBIF in 2 are sometimes skewed due to the pagination method that the networks use to supply users with their dataset registries. Registry pages contained set amounts (e.g. 20) of URLs and represent small slices of the actual data network registry. For registries that use pagination, the observatory would keep querying for registry pages until reaching the page or failing to respond. For instance, GBIF's URL and dataset totals in March 2019 (2.c) are low because an early query to a GBIF registry page was not answered and, consequently, the URLs of registry pages that should have followed were not discovered. Similar events happened for both the GBIF and DataONE observatories at later points in time, potentially overestimating the reliability of the data network.

In an effort to minimize artificial link rot due to internet access issues in our local network, we deployed the Preston observatories in a large commercial data center in Germany.

# Discussion <sub>338</sub>

We have shown that the reliability of URLs decreases over time in all of <sub>339</sub> the major biodiversity data networks that we monitored. If current trends <sub>340</sub> continue, the extent of reference rot will only worsen. Systematic changes <sub>341</sub> in the way we preserve and reference data are needed to reverse these <sub>342</sub> trends and improve the longevity and long-term integrity of the <sub>343</sub> biodiversity data record. Before we propose such changes, it's necessary to <sub>344</sub> first understand why URLs are proving to be ill-suited for referencing data <sub>345</sub> in the long term. <sub>346</sub>

## Unreliability of Location-based Identifiers <sub>347</sub>

The problems related to using URLs for referencing datasets are largely <sub>348</sub> due to the fact that they are location-based identifiers; they describe where <sub>349</sub> the data is but not necessarily what it is. Also, by definition, data accessed <sub>350</sub> via URLs must be mediated by a central authority, such as the <sub>351</sub> institutional repositories that serve biodiversity datasets, who can match <sub>352</sub> location-based identifiers with data. Interested users are expected to trust <sub>353</sub> the central authority to guarantee long-term access to the referenced data <sub>354</sub> in its original form. <sub>355</sub>

The use of URLs as identifiers violates the requirements of uniqueness <sub>356</sub> and persistence [Paskin, 1999]. An identifier must only ever identify one <sub>357</sub> entity (uniqueness) and must persist longer than the entity it identifies <sub>358</sub> (persistence) [Paskin, 1999]. However, as we have shown in our <sub>359</sub> experiments, many URLs do not possess both uniqueness and persistence; <sub>360</sub> unstable URLs forfeit uniqueness in the event of content drift, while <sub>361</sub> unresponsive URLs do not persist as long as the datasets they identify. <sub>362</sub>

15

At the core of URL instability is the current practice of using URLs to identify evolving datasets rather than fixed dataset versions. If biodiversity data providers were uniformly committed to allocating one URL per dataset version, then content drift might indeed become far less common, improving overall URL stability; however, widespread social adoption of such a commitment from all data providers may be unrealistic. Additionally, such a commitment would not address link rot and URL unresponsiveness. Even if a similar commitment were made by data providers to guarantee the long-term responsiveness of URLs, it could not address the case where a data provider either loses authority over a domain name or migrates to another. For example, our deployed Preston observatories recorded the sudden migration of over 14,000 Plazi datasets from the http://plazi.cs.umb.edu/ domain to http://tb.plazi.org/, an event which invalidated any references to URLs within the first domain.

Paskin proposed that "the best way to 'future proof' an identifier scheme is to forego any intelligence within the identifier itself" [Paskin, 1999], where the notion of intelligence refers to the inclusion of meaningful information in the textual representation of the identifier. URLs are structured according to the Domain Name System specification and inherently contain some minimum amount of intelligence: the domain that the URL belongs to [Mockapetris, 1987]. Thus, it is necessary to look to another identification scheme to allow for proper identification and reliable referencing.

An Alternative: Unique Content-Based Identifiers Instead of identifying digital datasets by location (i.e. URL), we can identify datasets by their content. One way to achieve this is to use algorithmically generated

content-based identifiers. A variety of cryptographic hashing algorithms 389
are available which guarantee a single unique hash, representable as text, 390
for any given dataset [NIST, 2001]. Because the hash itself is 391
deterministically derived from the content it identifies, we say that it is a 392
content-based identifier. Because hashes are deterministic, anyone 393
interested in identifying a dataset can simply compute its hash without the 394
need for some mediating central authority [Paskin, 1999]. If a change is 395
made to the dataset, then the hash computed from the modified dataset 396
will be different from that of the original. Therefore, if the hash of a 397
dataset is the same as the referenced hash, it must be the originally 398
referenced dataset [NIST, 2001]. Because hash identifiers can only identify 399
the exact content that was referenced, content drift is impossible; a content 400
hash will never match with either a different version of the content any 401
other content. Additionally, the chance of link rot is diminished due to the 402
lack of a single point of failure in the form of a central authority that is 403
solely responsible for making content available. The shift from 404
location-based to content-based identifiers allows for the decoupling of 405
future dataset accessibility from the original point of access. As long as 406
there exists some well-known and accessible data repository that has 407
archived the desired content, it can always be retrieved. Even if one 408
repository becomes inaccessible, another may be available to retrieve the 409
content. If a repository changes location, the reference is still reliable; it is 410
the interested user's responsibility to find either the repository's new 411
location or another repository that hosts the desired dataset. Additionally, 412
it is worth noting that duplication of content across different information 413
platforms does not lead to ambiguous references, but rather to distributed 414

17

copies of the same reliably addressed content. Figure 4 demonstrates the <sub>415</sub>

differences in referenced dataset retrieval when using location- versus <sub>416</sub>

content-based identifiers. <sub>417</sub>

(insert figure 4, see appendix) <sub>418</sub>

## Transitioning to Reliable References <sub>419</sub>

Although we propose a change in the fundamental mechanisms used to <sub>420</sub>

reference datasets, existing references can be made reliable with only minor <sub>421</sub>

modifications. Consider the following citation generated by GBIF <sub>422</sub>

according to their citation guidelines [GBIF.org, 2019]: <sub>423</sub>

Levatich T, Padilla F (2017). EOD - eBird Observation <sub>424</sub>

Dataset. Cornell Lab of Ornithology. Occurrence dataset <sub>425</sub>

https://doi.org/10.15468/aomfnb accessed via GBIF.org on <sub>426</sub>

2018-09-02. <sub>427</sub>

The citation references the eBird dataset hosted at gbif.org as it was <sub>428</sub>

retrieved on September 11, 2018. However, at the time of writing, the URL <sub>429</sub>

https://doi.org/10.15468/aomfnb redirects to a GBIF internal reference <sub>430</sub>

page which states that the eBird dataset was last updated in March of <sub>431</sub>

2019. The dataset made available through the listed URL is different from <sub>432</sub>

what was originally referenced in the citation, but it is impossible to <sub>433</sub>

determine the extent of the changes without having access to previous <sub>434</sub>

versions of the data. <sub>435</sub>

Fortunately, references like the example above can be made more <sub>436</sub>

reliable by augmenting them with a content-based identifier for the dataset. <sub>437</sub>

Consider the following enriched citation for the eBirds dataset adds a <sub>438</sub>

SHA-256 content hash [NIST, 2001]: <sub>439</sub>

18

> Levatich T, Padilla F (2017). EOD - eBird Observation
> Dataset. Cornell Lab of Ornithology. Occurrence dataset
> hash://sha256/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c
> accessed at https://doi.org/10.15468/aomfnb via GBIF.org on
> 2018-09-02.

440
441
442
443
444

The content hash is captured in a content address URI in the form of    445
hash://algo/hash-string proposed by [Trask, 2015], where "algo" is a    446
hashing algorithm (e.g., "sha256") and "hash-string" is the content hash    447
generated by the algorithm. In the example above, the hashing algorithm    448
is SHA256 and the hash string starts with 29d3. The added content hash    449
was derived from and uniquely identifies the exact version of the eBird    450
dataset that was originally referenced. If an interested user knows of and    451
has access to an information retrieval system that has indexed the dataset,    452
finding the desired dataset is as simple as querying for its content hash.    453
With the addition of a content hash, the URL becomes superfluous and is    454
included merely to demonstrate that the URL and content hash are not    455
mutually exclusive.    456

## Enhancing Dataset References with Provenance    457

A dataset reference can be given enhanced context by also referencing the    458
record that describes its provenance. The following citation further    459
augments the eBird dataset reference with the content hash of an    460
associated provenance record:    461

> Levatich T, Padilla F (2017). EOD - eBird Observation    462
> Dataset. Cornell Lab of Ornithology. Occurrence dataset    463

hash://sha256/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c accessed at https://doi.org/10.15468/aomfnb via GBIF.org on 2018-09-02 with provenance hash://sha256/b83cf099449dae3f633af618b19d05013953e7a1d7d97bc5ac01afd7bd9abe5d.

As was the case for the dataset, the provenance itself can be retrieved by querying a well-known information system that has indexed the hash of the referenced provenance record. Note that the provenance hash is not strictly necessary to make a dataset reference reliable; the dataset hash alone is sufficient. However, explicitly referencing the provenance of the dataset is useful because it allows future readers to also retrieve the same context that the original researcher who referenced the dataset had access to. More generally, the provenance describes the context of the retrieval of any type of content (e.g. datasets, metadata, citation files, etc.). The types of information in the provenance depend on the implementation of the data observatory, but at a minimum include the URLs that were queried to produce the content, the dates of the queries, the format of the content, and the data registries that were searched to find the content.

(insert figure 5, see appendix)

The use cases for the included provenance hash are many. For example, if the provenance record of a dataset is found, it may be possible to traverse the provenance and find newer versions of the dataset. This requires that the various versions of the dataset were observed at some point in time by a provenance-generating data observatory, properly archived, then made publicly accessible.

Our proposal to use Trask's content-addressed URIs to reliably reference data is similar to, and was inspired by, Kuhn & Dumontier's

20

method to make digital content verifiable and permanent using trusty <sub>490</sub> URIs [Kuhn and Dumontier, 2015]. We chose to use Trask's content hash <sub>491</sub> URIs because they are location and content agnostic and easy to read. <sub>492</sub> However, we recognize that trusty URIs can help facilitate content <sub>493</sub> retrieval and processing using a location-based URI prefix and an <sub>494</sub> (optional) extension suffix respectively. <sub>495</sub>

## Dataset Retrieval Using Hash References <sub>496</sub>

The dataset and provenance hashes referenced in the sample references <sub>497</sub> above were produced by our Preston observatories which were set up to <sub>498</sub> monitor the four data networks. Both the referenced dataset and its <sub>499</sub> provenance are available online at zenodo.org [Poelen, 2019f, Poelen, <sub>500</sub> 2019e, Poelen, 2019c] and archive.org [Poelen, 2019d]. A query for the <sub>501</sub> provenance hash in the search bar at zenodo.org or hash-archive.org should <sub>502</sub> direct the user to an archived repository of Preston observations that <sub>503</sub> contains both the dataset and its provenance (5). Given Zenodo's <sub>504</sub> long-term guarantee for data persistence and version availability [Zenodo, <sub>505</sub> 2019], the dataset reference is now reliable; it is effectively immune to both <sub>506</sub> link rot and content drift. Future readers can trust that the dataset will <sub>507</sub> stay available and, when downloaded, identically match the exact version <sub>508</sub> of the eBird dataset we referenced. Note that, to comply with Zenodo's <sub>509</sub> limitations on user uploads [Zenodo, 2019], we only exposed the set of <sub>510</sub> provenance hashes collected by each deployed Preston observatory for <sub>511</sub> search indexing, which are far fewer in number than the dataset hashes. <sub>512</sub> Thus, a query to zenodo.org for the dataset hash above should not produce <sub>513</sub> any results. This is an artificial limitation; ideally, an information system <sub>514</sub>

would index the dataset hashes as well. Note that our Zenodo publication 515
for the GBIF/iDigBio/BioCASe observatory [Poelen, 2019c] contains only 516
provenance, although the Internet Archive publication [Poelen, 2019d] 517
contains the content as well as provenance. Our Zenodo and Internet 518
Archive publications for BHL [Poelen, 2019e, Poelen, 2019a] and 519
DataONE [Poelen, 2019f, Poelen, 2019b] contain both content and 520
provenance. 521

Several biodiversity data aggregators, such as GBIF and iDigBio, 522
produce a citation file for each user query to allow researchers to simply 523
reference a single citation file rather than each individual dataset. A 524
citation file lists the URLs of the datasets (among other things, such as 525
attributions and retrieval dates) that were retrieved by the issued query. 526
We have demonstrated that dataset URLs are unreliable references; thus, 527
citation files that rely on URLs as references are also unreliable. Citation 528
files could be made reliable if they were augmented with the hashes of the 529
retrieved datasets and, optionally, their provenance records. In fact, 530
citation files themselves can be referenced by hash, along with 531
accompanying provenance hashes, as long as they are archived and made 532
accessible. 533

## DOIs for Datasets and Queries 534

Biodiversity data aggregators often assign each dataset or query a Digital 535
Object Identifier (DOI) [Paskin, 2009] (e.g. 10.123/456) wrapped as URL 536
(e.g. https://doi.org/10.123/456) and advise researchers to reference the 537
generated DOI rather than a URL. Unfortunately, this abstraction does 538
little to enhance the reliability of the reference. 539

22

The DOI Handle System [Paskin, 2009] associates DOIs with online    540
resources. However, it does not enforce any constraint on type of resource    541
associated with a DOI. When DOIs are used to reference biodiversity    542
datasets, the associated resources are often URLs, and therefore the use of    543
such DOIs as referencing mechanisms is just as potentially unreliable as    544
using URLs. In practice, these DOIs identify the evolving dataset (or set    545
of datasets in the case of a query) rather than a fixed version, as    546
demonstrated in the example references above. It is possible that an    547
author would wish to make such a reference to an evolving online digital    548
object. For example, an author promoting use of a published dataset might    549
want future users to be directed to the most up-to-date content. However,    550
such a fluid reference is not appropriate for making published results    551
reproducible.    552

The Handle System allows for a complex web of redirection and    553
distributed responsibilities. Just as the Domain Name System resolves    554
URLs to IP addresses, the Handle System resolves DOIs to data. When    555
these data are URLs, they must then be resolved through the Domain    556
Name System in order to retrieve the referenced content. However, the    557
responsibility for resolving DOIs to URLs is divided between the Handle    558
System and DOI registrars. The Handle System serves as the central    559
authority that maps DOI prefixes to DOI registrars, examples of which    560
include BHL, DataONE, GBIF, and iDigBio. These registrars are then    561
responsible, and indeed the central authorities for, associating DOIs that    562
match their designated prefix with URLs, and are free to change the URL    563
associated with any given DOI under their jurisdiction [Paskin,    564
2009, Foundation, 2018].    565

23

The ability of biodiversity data networks to change the URL associated <sub>566</sub> with a DOI is good for reference reliability in the sense that networks can <sub>567</sub> account for dataset migration without compromising existing references. <sub>568</sub> However, the use of DOIs addresses neither the instability of the URLs <sub>569</sub> they redirect to nor cases of link rot in which no URLs remain responsive <sub>570</sub> to serve the referenced dataset. Additionally, as the number of datasets <sub>571</sub> identified online continues to grow, proper maintenance of all of the DOIs <sub>572</sub> a data network administrates might become more unsustainable over time, <sub>573</sub> potentially increasing the risk of unreliable URLs going undetected. <sub>574</sub>

In an article proposing HTTP-URI-based stable identifiers (e.g. URLs <sub>575</sub> that are resolvable over HTTP) for biological collection objects, Güntsch et <sub>576</sub> al. admit that the use of DOIs does not solve the problem of unreliable <sub>577</sub> referencing but merely deflects the burden of URL maintenance onto <sub>578</sub> institutional repositories [Güntsch et al., 2017]. In contrast, we propose a <sub>579</sub> dataset referencing scheme that is reliable and can be supported by existing <sub>580</sub> infrastructures and workflows. If existing workflows require references to <sub>581</sub> be in the form of DOIs, it could be convenient to embed content hashes <sub>582</sub> into DOIs. Such an approach has already been established for ISBNs <sub>583</sub> through the creation of actionable ISBNs, or ISBN-As [Weissberg, 2008], <sub>584</sub> which may serve as a model for actionable content hashes. <sub>585</sub>

## What it Means to Preserve Data <sub>586</sub>

Our results indicate that reference rot poses an existential threat to <sub>587</sub> published biodiversity datasets. We've seen that the use of content-based <sub>588</sub> identifiers can effectively address the issue of reference rot. However, <sub>589</sub> identifiers are of little use in a vacuum. An identifier can only be useful for <sub>590</sub>

data retrieval when combined with a resolver to associate identifiers with 591 locations and a database to retrieve the dataset at the associated 592 location [Paskin, 1999]. Thus, we need to address how resolvers and 593 databases might be organized to accommodate content-based identifiers in 594 order to fully realize long-term data preservation. In this context, we 595 define data preservation as the continued capacity for datasets to be 596 reliably referenced and retrieved in their original form even as the global 597 digital biodiversity network evolves over time. 598

We propose four requirements that must be met to ensure proper data 599 preservation that prevents data loss: 1) datasets must be addressable and 600 retrievable using content-based rather than location-based identifiers; 2) an 601 agent must exist to collect datasets, record their provenance, and deposit 602 both to a dedicated repository; 3) these repositories should archive data 603 rather than discarding it; and 4) well-known search indexes should be 604 available to resolve hash identifiers to dataset locations within such 605 repositories. For the purposes of archival, it is important that the recorded 606 provenance records do not necessarily describe the datasets themselves, but 607 rather the activities that led to the procurement of those datasets; the 608 primary purpose of provenance in the context of an archive is to document 609 the fact that evidence, i.e. the dataset itself, does exist and to make it 610 discoverable for interested users [Bearman, 1995]. 611

We have shown that software agents such as Preston can be used to 612 collect datasets and their provenance over time while maintaining 613 content-addressability; all that is needed to ensure proper data 614 preservation are a dedicated repository and a well-known, publicly 615 available search index to map content-based identifiers to datasets located 616

25

in the repository. In practice, repositories and search indexes (and 617
potentially software agents such as Preston deployments) can be 618
co-located; examples include Zenodo and the Internet Archive, although 619
they impose some limitations that may restrict file size, number of files, 620
and the amount of information that can be indexed [Zenodo, 2019, Archive, 621
2019]. These existing information systems may serve as models for 622
long-term biodiversity information systems. 623

These requirements help to ensure that biodiversity data remain FAIR 624
(Findable, Accessible, Interoperable, and Reusable) [Wilkinson et al., 2016]. 625
Findability is achieved through the publishing of provenance logs which 626
thoroughly describe what datasets are and where they originated from. 627
The amenability of the content-based identification paradigm to the 628
operation of independent distributed repositories strengthens accessibility 629
by preventing the failure of a single data repository from inhibiting future 630
data access (4). Content-based identification also allows for interoperability 631
due to the absence of any central authority to administrate data access; a 632
content hash computed from a dataset is guaranteed to match the hash 633
computed by any other agent using the same dataset. Finally, and 634
particularly relevant to this paper's purpose, reusability is strengthened by 635
enhancing the retrievability of referenced datasets and allowing users to 636
verify that a retrieved dataset exactly matches that which was referenced. 637

# Conclusions 638

Although reference rot is resulting in a steady decline in the reliability of 639
our digital biodiversity record, realistic solutions are available to address 640
the root causes of the issue. Content drift can be eliminated altogether by 641

26

changing the way we reference datasets, from using location-based 642 identifiers to ones that are content-based. Meanwhile, the online 643 biodiversity data networks can be made far more resilient to link rot if 644 distributed observation and archival techniques are used to capture 645 incremental changes to the data record so that references can remain valid 646 even when online datasets are updated, removed, or relocated. 647

The use of content-based identifiers should be considered by biodiversity 648 data aggregators in order to increase the reliability of references to the 649 data they aggregate. If long-term data observatories for biodiversity data 650 networks are established, their collected data routinely deposited to 651 well-known publicly available archives, and the archived data sufficiently 652 indexed, then researchers and data curators will be able to have certainty 653 that the datasets they contribute and reference will maintain reliability in 654 the midst of an ever-changing digital ecosystem. 655

Great care has been taken to establish rigorous preservation guidelines 656 for physical specimens, yet there is much that can be done to increase the 657 longevity of our digital data. Our method is not only suited for tracking 658 datasets in biodiversity data networks, but also provides a resilient and 659 reliable way to publish, reference, and preserve scientific digital datasets 660 without having to abandon our existing infrastructures. The method 661 provides a much-needed foundation for constructing digital provenance 662 graphs from an accessible, verifiable, and citable digital scholarly record. 663

# Acknowledgments 664

27

# References

Archive, 2019. Archive, I. (2019). Uploading - a basic guide. Accessed: 2019-12-04.

Bearman, 1995. Bearman, D. (1995). Archival strategies. *The American Archivist*, 58(4):380–413.

Berners-Lee et al., 2005. Berners-Lee, T., Fielding, R. T., and Masinter, L. M. (2005). Uniform Resource Identifier (URI): Generic Syntax. RFC 3986.

Berners-Lee et al., 1994. Berners-Lee, T., Masinter, L. M., and McCahill, M. P. (1994). Uniform Resource Locators (URL). RFC 1738.

Costello et al., 2013. Costello, M. J., Bouchet, P., Boxshall, G., Fauchald, K., Gordon, D., Hoeksema, B. W., Poore, G. C. B., van Soest, R. W. M., Stöhr, S., Walter, T. C., Vanhoorne, B., Decock, W., and Appeltans, W. (2013). Global coordination and standardisation in marine biodiversity through the world register of marine species (WoRMS) and related databases. *PLoS ONE*, 8(1):e51629.

DataONE, 2012. DataONE (2012). Dataone citation guidelines. Accessed: 2019-12-04.

Davis and Schmidt, 1996. Davis, E. B. and Schmidt, D. (1996). *Guide to Information Sources in the Botanical Sciences*. Vol. 2nd ed. Reference Sources in Science and Technology. Englewood, Colo: Libraries Unlimited.

Edwards, 2000. Edwards, J. L. (2000). Interoperability of biodiversity databases: Biodiversity information on every desktop. *Science*, 289(5488):2312–2314.

Elton, 1927. Elton, C. S. (1927). *Animal ecology.* Macmillan Co.,.

Facility, 2019. Facility, G. T. G. B. I. (2019). What is gbif? Accessed: 2019-12-04.

Foundation, 2018. Foundation, I. D. (2018). Doi handbook. Accessed: 2019-12-04.

Garfield et al., 1964. Garfield, E., Sher, I. H., and Torpie, R. J. (1964). *The Use of Citation Data in Writing the History of Science.* Institute for Scientific Information Inc Philadelphia PA.

GBIF.org, 2019. GBIF.org (2019). Gbif citation guidelines. Accessed: 2019-12-04.

Güntsch et al., 2017. Güntsch, A., Hyam, R., Hagedorn, G., Chagnoux, S., Röpert, D., Casino, A., Droege, G., Glöckler, F., Gödderz, K., Groom, Q., Hoffmann, J., Holleman, A., Kempa, M., Koivula, H., Marhold, K., Nicolson, N., Smith, V. S., and Triebel, D. (2017). Actionable, long-term stable and semantic web compatible identifiers for access to biological collection objects. *Database*, 2017.

Hortal et al., 2015. Hortal, J., de Bello, F., Diniz-Filho, J. A. F., Lewinsohn, T. M., Lobo, J. M., and Ladle, R. J. (2015). Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual Review of Ecology, Evolution, and Systematics*, 46(1):523–549.

iDigBio.org, 2016. iDigBio.org (2016). idigbio citation guidelines. Accessed: 2019-12-04.

Klein et al., 2014. Klein, M., de Sompel, H. V., Sanderson, R., Shankar, H., Balakireva, L., Zhou, K., and Tobin, R. (2014). Scholarly context not found: One in five articles suffers from reference rot. *PLoS ONE*, 9(12):e115253.

Kuhn and Dumontier, 2015. Kuhn, T. and Dumontier, M. (2015). Making digital artifacts on the web verifiable and reliable. *IEEE Transactions on Knowledge and Data Engineering*, 27(9):2390–2400.

Matsunaga et al., 2013. Matsunaga, A., Figueiredo, R., Thompson, A., Traub, G., Beaman, R., and Fortes, J. A. (2013). Integrated digitized biocollections (idigbio) cyberinfrastructure status and futures. In *TDWG 2013 ANNUAL CONFERENCE*.

Michener et al., 2011. Michener, W., Vieglais, D., Vision, T., Kunze, J., Cruse, P., and Janée, G. (2011). DataONE: Data observation network for earth: Preserving data and enabling innovation in the biological and environmental sciences. *D-Lib Magazine*, 17(1/2).

Mockapetris, 1987. Mockapetris, P. (1987). Domain names - concepts and facilities. RFC 1034.

NIST, 2001. NIST (2001). Descriptions of sha-256, sha-384, and sha-512. Accessed: 2019-12-04.

Page et al., 2015. Page, L. M., MacFadden, B. J., Fortes, J. A., Soltis, P. S., and Riccardi, G. (2015). Digitization of biodiversity collections reveals biggest data on biodiversity. *BioScience*, 65(9):841–842.

Paskin, 1999. Paskin, N. (1999). Toward unique identifiers. *Proceedings of the IEEE*, 87(7):1208–1227.

Paskin, 2009. Paskin, N. (2009). Digital object identifier (DOI®) system. In *Encyclopedia of Library and Information Sciences, Third Edition*, pages 1586–1592. CRC Press.

Poelen et al., 2018. Poelen, J., Elliott, M., Alzuru, I., and Patel, P. (2018). Preston: a biodiversity dataset tracker.

Poelen, 2019a. Poelen, J. H. (2019a). A biodiversity dataset graph: Biodiversity Heritage Library (BHL).

Poelen, 2019b. Poelen, J. H. (2019b). A biodiversity dataset graph: DataONE.

Poelen, 2019c. Poelen, J. H. (2019c). A biodiversity dataset graph: GBIF, iDigBio, BioCASe.

Poelen, 2019d. Poelen, J. H. (2019d). Biodiversity Dataset Archive.

Poelen, 2019e. Poelen, J. H. (2019e). A biodiversity dataset graph: Bhl.

Poelen, 2019f. Poelen, J. H. (2019f). A biodiversity dataset graph: Dataone.

Poelen, 2019g. Poelen, J. H. (2019g). To connect is to preserve: on frugal data integration and preservation solutions.

Rinaldo and Norton, 2009. Rinaldo, C. and Norton, C. (2009). BHL, the biodiversity heritage library: An expanding international collaboration. *Nature Precedings*.

Robertson et al., 2014. Robertson, T., Döring, M., Guralnick, R., Bloom, D., Wieczorek, J., Braak, K., Otegui, J., Russell, L., and Desmet, P. (2014). The GBIF integrated publishing toolkit: Facilitating the efficient publishing of biodiversity data on the internet. *PLoS ONE*, 9(8):e102623.

Trask, 2015. Trask, B. (2015). Principles of content addressing. `https://bentrask.com/?q=hash://sha256/` `98493caa8b37eaa26343bbf73f232597a3ccda20498563327a4c3713821df892`. Accessed: 2019-12-04.

Vision, 2010. Vision, T. J. (2010). Open data and the social contract of scientific publishing. *BioScience*, 60(5):330–331.

Weissberg, 2008. Weissberg, A. (2008). The identification of digital book content. *Publishing Research Quarterly*, 24(4):255–260.

Wieczorek et al., 2012. Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., Giovanni, R., Robertson, T., and Vieglais, D. (2012). Darwin core: An evolving community-developed biodiversity data standard. *PLoS ONE*, 7(1):e29715.

Wilkinson et al., 2016. Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., Bouwman, J., Brookes, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T., Finkers, R., Gonzalez-Beltran, A., Gray, A. J., Groth, P., Goble, C., Grethe, J. S., Heringa, J., 't Hoen, P. A., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S. J., Martone, M. E., Mons, A., Packer, A. L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone, S.-A.,
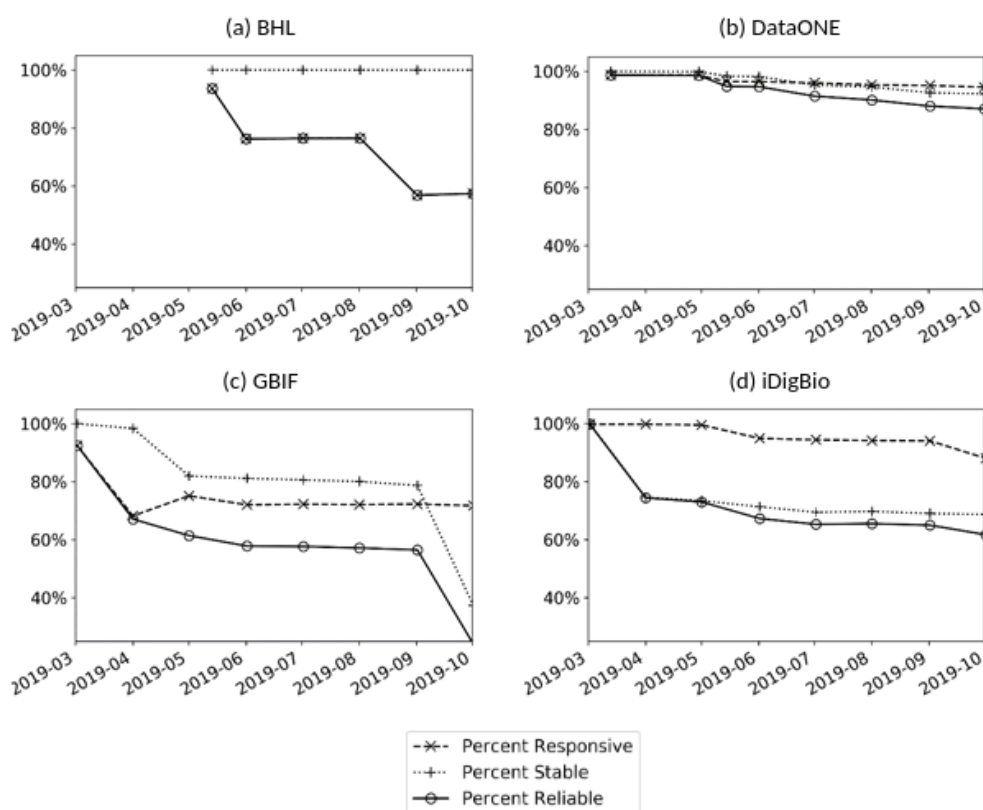
Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M. A., Thompson, M., van der Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., and Mons, B. (2016). The FAIR guiding principles for scientific data management and stewardship. *Scientific Data*, 3(1).

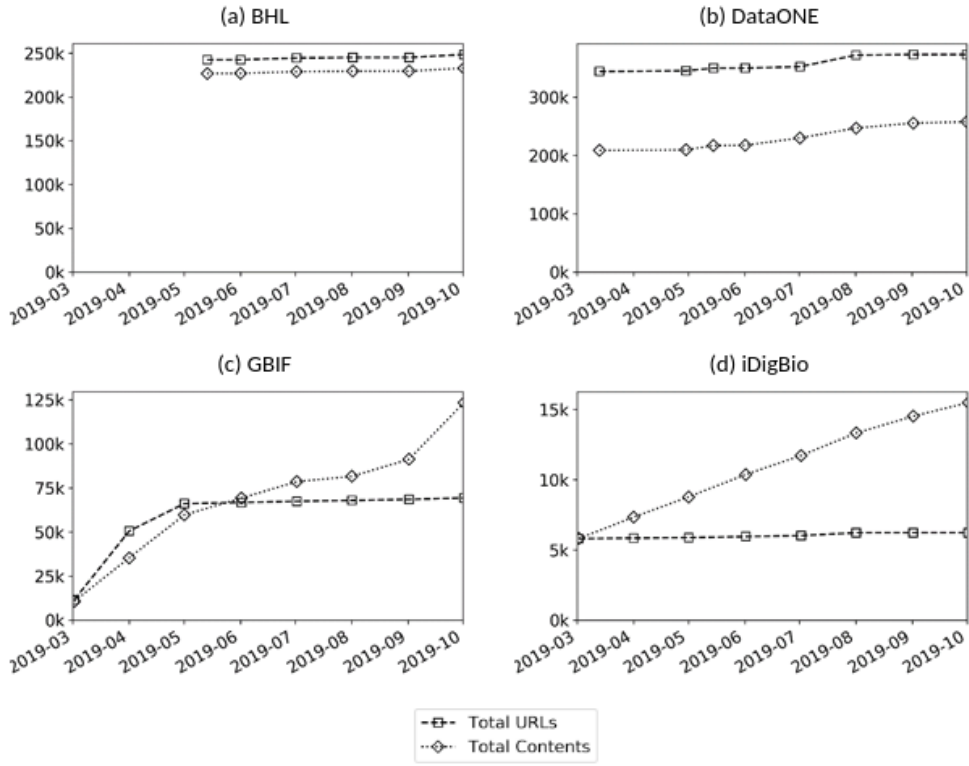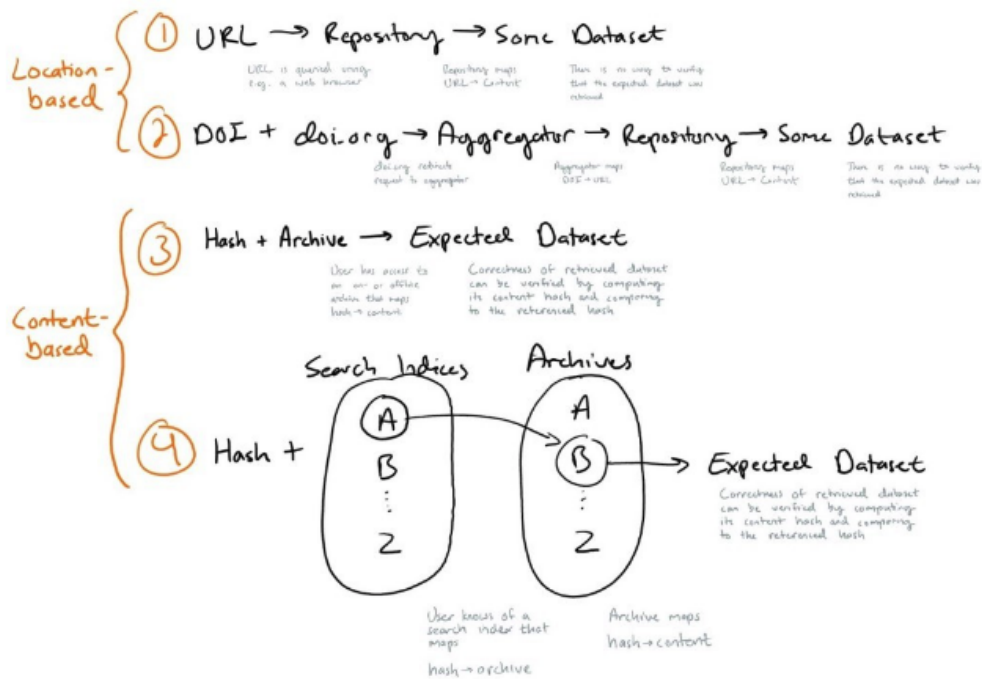Zenodo, 2019. Zenodo (2019). General policies. Accessed: 2019-12-04.

# Figures



**Figure 1.** Reliable record keeping for digital datasets (b) can be achieved in an analogous way to current practices in record keeping for physical specimens (a). Biologists collect physical specimens from the natural world, thoroughly document the process, then store the specimens in facilities equipped for long-term preservation. Analogously, digital datasets that are downloaded from the internet can be thoroughly documented and archived in dedicated repositories for long-term preservation. Just as the collection of physical specimens is recorded and identified in specimen history records, the downloading of digital datasets can also be recorded and identified in dataset history records.

**Figure 2.** Overall responsiveness, stability, and reliability from March 2019 to October 2019 as a percentage of URLs that exhibit each indicator in a) BHL, b) DataONE, c) GBIF, and d) iDigBio.

**Figure 3.** Total number of URLs and unique contents observed from March 2019 to October 2019 for a) BHL, b) DataONE, c) GBIF, and d) iDigBio.

**Figure 4.** Visualization of content resolution for location- versus content-based identifiers. 1) URLs point to a known location of a dataset, but do not guarantee either the presence or authenticity of the retrieved dataset; 2) the use of a DOI that resolves to a URL adds a layer of redirection; 3) A content-addressed dataset can be found by matching against recomputed hashes of available datasets in an archive; 4) well-known (online) hash indices can be used to facilitate discovery of dataset locations associated with a specific content hash.

**Figure 5.** An example of a search index mapping hashes to archives. A search for a content or provenance hash at hash-archive.org will find any associated URLs that have been registered at hash-archive.org.

# Tables

| Data Network | Responsive URLs | Stable URLs* | Reliable URLs |
|---|---|---|---|
| BHL | 57.41% (142,672) | 99.97% (232,996) | 57.39% (142,633) |
| DataONE | 94.55% (352,438) | 92.27% (339,109) | 87.09% (324,641) |
| GBIF | 71.72% (49,707) | 37.35% (20,094) | 24.05% (16,669) |
| iDigBio | 88.04% (5,477) | 68.69% (4,251) | 61.68% (3,837) |
| All observed URLs | 78.94% (546,645) | 90.43% (593,469) | 70.07% (485,203) |

**Table 1.** Overall responsiveness, stability, and reliability for URLs observed in each biodiversity data network and for all observed URLs as of October 2019. * URLs that never provided content were omitted from the divisor when calculating Stable URLs percentages.