

Toward Reliable Biodiversity Data References

Michael Elliott¹, Jorrit H. Poelen^{2,*}, José A.B. Fortes¹

**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.

Keywords— Biodiversity, Ecological Informatics, Information
Systems, Information Retrieval

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
([15], [7]). Records that describe these carefully collected specimens are
now made available digitally through online search indices, registries, and
data archives ([23]). 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” ([9]).
So far, various initiatives have succeeded to provide comprehensive
aggregate views from previously scattered natural history record siloes
([34], [20], [8], [19], [10]). 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 ([39], FORCE11, [35]) and periodic archival ([5]),
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
([12]). 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 ([13], [16], [6]), 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
online data networks for the matching author and title, but there is no
guarantee that information found this way will be exactly the same as

what was originally referenced. Any reference that does not allow future
readers to find the referenced work fails to satisfy the FAIR principle of
findability: "F1. (meta)data are assigned a globally unique and eternally
persistent identifier." ([40]). Our study is not alone in providing evidence
that suggests that networked, location-based access to digital objects is an
unreliable mechanism for providing continued access to the unaltered
original work ([37], [17]). Unless we change the way we preserve and cite
our digital scholarly works, our physical records stored in libraries and
museums around the world are likely to outlast our digital ones.

Problem Characterization

We show that the current practice of using Uniform Resource Locators
(URLs) ([4]) to reference online biodiversity datasets provides no
guarantee of long-term data accessibility. Readers who encounter
references that use URLs as dataset identifiers cannot be certain that the
referenced data will continue to be accessible and in its exact original form.
This uncertainty might be cause for alarm for researchers because, over
time, the integrity of the scholarly record itself is damaged when existing
references become reliable due to the loss of access to the data they
reference. When data access is lost, it is possible that documented research
results may become impossible to reproduce and the justification for any
conclusions or hypotheses that relied on lost results may be undermined. If
the use of error-prone referencing techniques is not addressed, we expect
that any resulting gaps in the biodiversity data record will only become
more severe.

The current practice of relying on URLs to locate and identify

referenced data is hazardous due to their demonstrated risk of link rot and
content drift ([17]). Link rot occurs when a URL, or link, that had
previously responded to queries can no longer be reached. This can
happen, for example, due to temporary outages, URL retirement, or URL
migration. A link exhibits content drift when a query to the link provides
content that is different from the content it provided in the past. The
extent of content drift can vary; content may have received only minor
edits with no changes in semantics, or it may reference a different entity
altogether. When a single URL is used to locate data that may change
over time, access to any particular version of the data is likely to be
short-lived. We show that, in the event of link rot or content drift, any
existing references that relied affected URL may become unreliable.

In one study on the Genetics journal, it was reported that 40% of links
(URLs) to supplemental materials became unavailable due to link rot
within one year of publication ([37]). Another study ([17]) confirmed that
as many as one in five articles in journal of Science, Technology, and
Medicine provide references that exhibit either link rot and content drift
and refer to the existence of either as “reference rot”. Since existing
biodiversity references largely rely on URLs to locate datasets, it is
reasonable to expect that biodiversity data networks are also at risk of
providing unreliable dataset references as a result of reference rot. The
information systems used by major biodiversity data networks, such as
DataONE, GBIF, and iDigBio, rely on data curators, such as institutional
repositories, to maintain active dataset URLs, and aggregate the data
found at those URLs for distribution in response to user queries. If a data
curator modifies, relocates, or stops serving a particular dataset, it may

become impossible to retrieve the original dataset and the integrity of the
data network will suffer as a result.

In this paper, we propose a methodology for measuring the existence of
link rot and content drift in online data networks, then provide
experimental results that confirm the existence of both link rot and
content drift across all of the biodiversity data networks we considered,
including BHL, DataONE, iDigBio, and GBIF. Finally, we propose a
method for referencing and serving biodiversity data in a way that works
toward satisfying the Findable, Accessible, Interoperable, and Reusable
(FAIR) principles ([40]).

Methodology

Although it has been demonstrated that reference rot does occur when
URLs are used for referencing scientific works ([37], [17]), we are not
aware of any prior studies that provide quantitative evidence that reference
rot occurs specifically in biodiversity data networks. We set out to
quantify the extent of reference rot in biodiversity data networks. Because
reference rot occurs in the scope of individual data references, and
references to digital datasets rely on URLs to locate the data, we begin by
introducing terminology for characterizing the reliability of a URL
according to how often it exhibits link rot and content drift.

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 ([3]). 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

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

• Unstable: the content that the link points to sometimes changes 166

• Stable: the content that the link points to never changes 167

We characterize the overall reliability of a URL according to both of its 168
responsiveness and stability: 169

• Unreliable: the link does not always provide the expected content; it 170
is either unresponsive, unstable, or both 171

• Reliable: the link always provides the expected content; it is both 172
responsive and stable 173

Before we can determine the reliability of any given URL, we must first 174
monitor its behavior over time by documenting how it responds to periodic 175
queries. For the context of biodiversity, we consider the case when the 176
content that a URL produces is a dataset. 177

The Data Collection Process 178

We suggest that digital dataset collection practices have some analogies to 179
well-established physical specimen collection procedures (Fig. 1) ([33]). If 180
datasets are considered analogous to specimens, then the URLs that locate 181
datasets are analogous to the physical locations of specimens in the natural 182
world; they are where digital datasets were originally found, but not where 183
they should be preserved. Once found, physical specimens are collected by 184
hand; similarly, digital datasets are downloaded by querying their URLs. 185
Once a specimen is collected and deposited to a safe, well-known 186
repository, a record is kept that documents what the specimen is in 187
addition to when, where, and by whom it was collected. 188

(insert figure 1 / see appendix)

The same can be done for downloaded datasets. When a dataset is downloaded, a record can be kept that details the URL that was queried, the time of query, and who (e.g. a human or software agent) issued the query that initiated the download event; we refer to this record as the dataset’s provenance record. Additionally, the dataset itself should be stored in a safe, well-known dataset archive. The final step in the collection process is to link the actual preserved specimen to its corresponding record (the “specimen history” in Fig. 1) via an assigned unique identifier. For digital datasets, we use cryptographic hashes of the data as unique content-based identifiers.

Data Collection Over Time

By establishing a dedicated data observatory that follows the collection process we have described, we can build a history for each observed URL to capture its long-term reliability. Such an observatory should periodically query the URLs listed in data network’s URL registry, producing for each URL two complementary parts: 1) an archived copy of the response to the corresponding query, whether it was a dataset, an error code, or no reply at all, and 2) a record of its provenance, including the URL itself, the current date, and a content-based identifier of any dataset received. The use of a content-based data identifier is crucial; it allows us to reliably link each acquired dataset to its provenance record without the need for an intermediate index. Successive provenance records can be aggregated to construct comprehensive histories for both datasets (when and where they were found) and URLs (which datasets they produced over a series of

queries over time). 214

The constructed URL histories can be analyzed to determine whether a 215
link was ever broken, when it was broken, and whether it became 216
responsive again. The logs also identify the content (or lack of it) that a 217
URL produced each time it was queried. Any change in the content 218
identifier from query to the next indicates a change in the content of the 219
dataset. These link breakages and content changes correlate to link rot and 220
content drift, respectively, and allow us to determine the responsiveness, 221
stability, and reliability of each URL over time. 222

Data Network Reliability 223

Now that we have outlined a method for observing and documenting the 224
behavior of URLs over an extended period of time, we can apply our 225
method to observe all of URLs registered by biodiversity data networks. 226
We also extend the idea of URL reliability to entire data networks and 227
propose that the overall reliability of a data network can be evaluated by 228
monitoring the long-term reliability of each individual URL in the network 229
exposes. Whereas we rigidly label individual URLs with binary indicators 230
of responsiveness, stability, and reliability, we grade data networks 231
according to the percentage of registered URLs that are assigned each of 232
the reliability indicators. For example, if a data network contains three 233
distinct URLs and we find that only two out of the three are reliable, then 234
we say the data network is 67% reliable. 235

Experiment

236

The Preston biodiversity dataset tracker ([26]) implements mechanisms for monitoring data networks as we have described. It allows users to deploy a data network observatory which systematically observes the entire set of URLs registered by the network, queries each URL for data, then documents data collection and archives the results. All crawl activities, the queries they issue, and the results they produce are meticulously recorded in a string of provenance logs.

237

238

239

240

241

242

243

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

244

245

246

247

248

249

250

251

252

253

254

255

(insert figure 2, see appendix)

256

Results

257

Breakdowns of the overall reliabilities of the data networks are provided in Table 1. Results are listed as percentages and total counts of URLs in the

258

259

data network that were assigned each reliability indicator. When analyzing
the recorded results of queries to URLs in each data network over a period
of seven months, we found that, for each individual 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.

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

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

(insert figure 3, see appendix)

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, 286
 although BHL scored relatively low in responsiveness due to frequent link 287
 rot, the content that it does provide is more stable than all other networks 288
 because content drift within BHL is relatively rare. Conversely, although 289
 iDigBio is relatively responsive, it has low stability because the network's 290
 near-constant content growth far outpaces its URL growth. GBIF's 291
 behavior was characterized by large sporadic swings; a mass URL 292
 migration of over 14,000 Plazi-hosted datasets occurred in May, 293
 introducing thousands of new URLs over a short period of time, while over 294
 31,000 URLs (60% of URLs that responded to queries that month) 295
 suddenly changed contents in October. Even the most reliable network, 296
 DataONE, shows a clear downward trend in all three categories, with 13% 297
 of URLs becoming unreliable over a period of just seven months. 298
 Additionally, DataONE's growth curves indicate that there are far fewer 299
 unique contents than unique URLs; this evokes two possibilities: either 300
 much of DataONE's URL population is unresponsive, or DataONE lists 301
 multiple URLs for many of its datasets. Because DataONE has been 302
 shown to be highly responsive, it could be the case that many distinct 303
 URLs refer to the same datasets. It's also worth noting that the June and 304
 September spikes in BHL's unresponsiveness were largely due to URLs 305
 that failed to respond in those particular months but actually did respond 306
 to future queries. 307

Sources of Potential Numerical Error 308

We expect that the URL reliability counts generated for the figures and 309
 tables are lower than their actual values. When we qualified URLs as being 310

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

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

Unreliability of Location-based Identifiers

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

The use of URLs as identifiers violates the requirements of uniqueness and persistence ([24]). An identifier must only ever identify one entity (uniqueness) and must persist longer than the entity it identifies (persistence) ([24]). However, as we have shown in our experiments, many URLs do not possess both uniqueness and persistence; unstable URLs forfeit uniqueness in the event of content drift, while unresponsive URLs do not persist as long as the datasets they identify.

At the core of URL instability is the current practice of using URLs to 359
identify evolving datasets rather than fixed dataset versions. If biodiversity 360
data providers were uniformly committed to allocating one URL per 361
dataset version, then content drift might indeed become far less common, 362
improving overall URL stability; however, widespread social adoption of 363
such a commitment from all data providers may be unrealistic. 364
Additionally, such a commitment would not address link rot and URL 365
unresponsiveness. Even if a similar commitment were made by data 366
providers to guarantee the long-term responsiveness of URLs, it could not 367
address the case where a data provider either loses authority over a 368
domain name or migrates to another. For example, our deployed Preston 369
observatories recorded the sudden migration of over 14,000 Plazi datasets 370
from the <http://plazi.cs.umb.edu/> domain to <http://tb.plazi.org/>, an event 371
which invalidated any references to URLs within the first domain. 372

Paskin proposed that “the best way to ‘future proof’ an identifier 373
scheme is to forego any intelligence within the identifier itself” ([24]), 374
where the notion of intelligence refers to the inclusion of meaningful 375
information in the textual representation of the identifier. URLs are 376
structured according to the Domain Name System specification and 377
inherently contain some minimum amount of intelligence: the domain that 378
the URL belongs to ([21]). Thus, it is necessary to look to another 379
identification scheme to allow for proper identification and reliable 380
referencing. 381

An Alternative: Unique Content-Based Identifiers Instead of identifying 382
digital datasets by location (i.e. URL), we can identify datasets by their 383
content. One way to achieve this is to use algorithmically generated 384

content-based identifiers. A variety of cryptographic hashing algorithms
are available which guarantee a single unique hash, representable as text,
for any given dataset ([22]). Because the hash itself is deterministically
derived from the content it identifies, we say that it is a content-based
identifier. Because hashes are deterministic, anyone interested in
identifying a dataset can simply compute its hash without the need for
some mediating central authority ([24]). If a change is made to the
dataset, then the hash computed from the modified dataset will be
different from that of the original. Therefore, if the hash of a dataset is the
same as the referenced hash, it must be the originally referenced dataset
([22]). Because hash identifiers can only identify the exact content that
was referenced, content drift is impossible; a content hash will never match
with either a different version of the content any other content.
Additionally, the chance of link rot is diminished due to the lack of a single
point of failure in the form of a central authority that is solely responsible
for making content available. The shift from location-based to
content-based identifiers allows for the decoupling of future dataset
accessibility from the original point of access. As long as there exists some
well-known and accessible data repository that has archived the desired
content, it can always be retrieved. Even if one repository becomes
inaccessible, another may be available to retrieve the content. If a
repository changes location, the reference is still reliable; it is the
interested user's responsibility to find either the repository's new location
or another repository that hosts the desired dataset. Additionally, it is
worth noting that duplication of content across different information
platforms does not lead to ambiguous references, but rather to distributed

copies of the same reliably addressed content. Figure 4 demonstrates the
differences in referenced dataset retrieval when using location- versus
content-based identifiers.

(insert figure 4, see appendix)

Transitioning to Reliable References

Although we propose a change in the fundamental mechanisms used to
reference datasets, existing references can be made reliable with only minor
modifications. Consider the following citation generated by GBIF
according to their citation guidelines ([13]):

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

The citation references the eBird dataset hosted at gbif.org as it was
retrieved on September 11, 2018. However, at the time of writing, the URL
<https://doi.org/10.15468/aomfnb> redirects to a GBIF internal reference
page which states that the eBird dataset was last updated in March of
2019. The dataset made available through the listed URL is different from
what was originally referenced in the citation, but it is impossible to
determine the extent of the changes without having access to previous
versions of the data.

Fortunately, references like the example above can be made more
reliable by augmenting them with a content-based identifier for the dataset.
Consider the following enriched citation for the eBirds dataset adds a
SHA-256 content hash ([22]):

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/aomfmb via GBIF.org on
 2018-09-02.

The content hash is captured in a content address URI in the form of
 hash://algo/hash-string proposed by [36], where "algo" is a hashing
 algorithm (e.g., "sha256") and "hash-string" is the content hash generated
 by the algorithm. In the example above, the hashing algorithm is SHA256
 and the hash string starts with 29d3. The added content hash was derived
 from and uniquely identifies the exact version of the eBird dataset that
 was originally referenced. If an interested user knows of and has access to
 an information retrieval system that has indexed the dataset, finding the
 desired dataset is as simple as querying for its content hash. With the
 addition of a content hash, the URL becomes superfluous and is included
 merely to demonstrate that the URL and content hash are not mutually
 exclusive.

Enhancing Dataset References with Provenance

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

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/aomfmb 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

method to make digital content verifiable and permanent using trusty
 URIs ([18]). We chose to use Trask’s content hash URIs because they are
 location and content agnostic and easy to read. However, we recognize
 that trusty URIs can help facilitate content retrieval and processing using
 a location-based URI prefix and an (optional) extension suffix respectively.

Dataset Retrieval Using Hash References

The dataset and provenance hashes referenced in the sample references
 above were produced by our Preston observatories which were set up to
 monitor the four data networks. Both the referenced dataset and its
 provenance are available online at zenodo.org ([32], [31], [29]) and
 archive.org ([30]). A query for the provenance hash in the search bar at
 zenodo.org or hash-archive.org should direct the user to an archived
 repository of Preston observations that contains both the dataset and its
 provenance (5). Given Zenodo’s long-term guarantee for data persistence
 and version availability ([41]), the dataset reference is now reliable; it is
 effectively immune to both link rot and content drift. Future readers can
 trust that the dataset will stay available and, when downloaded, identically
 match the exact version of the eBird dataset we referenced. Note that, to
 comply with Zenodo’s limitations on user uploads ([41]), we only exposed
 the set of provenance hashes collected by each deployed Preston
 observatory for search indexing, which are far fewer in number than the
 dataset hashes. Thus, a query to zenodo.org for the dataset hash above
 should not produce any results. This is an artificial limitation; ideally, an
 information system would index the dataset hashes as well. Note that our
 Zenodo publication for the GBIF/iDigBio/BioCAsE observatory ([29])

contains only provenance, although the Internet Archive publication ([30]) 511
contains the content as well as provenance. Our Zenodo and Internet 512
Archive publications for BHL ([31], [27]) and DataONE ([32], [28]) 513
contain both content and provenance. 514

Several biodiversity data aggregators, such as GBIF and iDigBio, 515
produce a citation file for each user query to allow researchers to simply 516
reference a single citation file rather than each individual dataset. A 517
citation file lists the URLs of the datasets (among other things, such as 518
attributions and retrieval dates) that were retrieved by the issued query. 519
We have demonstrated that dataset URLs are unreliable references; thus, 520
citation files that rely on URLs as references are also unreliable. Citation 521
files could be made reliable if they were augmented with the hashes of the 522
retrieved datasets and, optionally, their provenance records. In fact, 523
citation files themselves can be referenced by hash, along with 524
accompanying provenance hashes, as long as they are archived and made 525
accessible. 526

DOIs for Datasets and Queries 527

Biodiversity data aggregators often assign each dataset or query a Digital 528
Object Identifier (DOI) ([25]) (e.g. 10.123/456) wrapped as URL (e.g. 529
<https://doi.org/10.123/456>) and advise researchers to reference the 530
generated DOI rather than a URL. Unfortunately, this abstraction does 531
little to enhance the reliability of the reference. 532

The DOI Handle System ([25]) associates DOIs with online resources. 533
However, it does not enforce any constraint on type of resource associated 534
with a DOI. When DOIs are used to reference biodiversity datasets, the 535

associated resources are often URLs, and therefore the use of such DOIs as 536
referencing mechanisms is just as potentially unreliable as using URLs. In 537
practice, these DOIs identify the evolving dataset (or set of datasets in the 538
case of a query) rather than a fixed version, as demonstrated in the 539
example references above. It is possible that an author would wish to make 540
such a reference to an evolving online digital object. For example, an 541
author promoting use of a published dataset might want future users to be 542
directed to the most up-to-date content. However, such a fluid reference is 543
not appropriate for making published results reproducible. 544

The Handle System allows for a complex web of redirection and 545
distributed responsibilities. Just as the Domain Name System resolves 546
URLs to IP addresses, the Handle System resolves DOIs to data. When 547
these data are URLs, they must then be resolved through the Domain 548
Name System in order to retrieve the referenced content. However, the 549
responsibility for resolving DOIs to URLs is divided between the Handle 550
System and DOI registrars. The Handle System serves as the central 551
authority that maps DOI prefixes to DOI registrars, examples of which 552
include BHL, DataONE, GBIF, and iDigBio. These registrars are then 553
responsible, and indeed the central authorities for, associating DOIs that 554
match their designated prefix with URLs, and are free to change the URL 555
associated with any given DOI under their jurisdiction. ([25], [11]) 556

The ability of biodiversity data networks to change the URL associated 557
with a DOI is good for reference reliability in the sense that networks can 558
account for dataset migration without compromising existing references. 559
However, the use of DOIs addresses neither the instability of the URLs 560
they redirect to nor cases of link rot in which no URLs remain responsive 561

to serve the referenced dataset. Additionally, as the number of datasets
identified online continues to grow, proper maintenance of all of the DOIs
a data network administrates might become more unsustainable over time,
potentially increasing the risk of unreliable URLs going undetected.

In an article proposing HTTP-URI-based stable identifiers (e.g. URLs
that are resolvable over HTTP) for biological collection objects, Güntsch et
al. admit that the use of DOIs does not solve the problem of unreliable
referencing but merely deflects the burden of URL maintenance onto
institutional repositories ([14]). In contrast, we propose a dataset
referencing scheme that is reliable and can be supported by existing
infrastructures and workflows. If existing workflows require references to
be in the form of DOIs, it could be convenient to embed content hashes
into DOIs. Such an approach has already been established for ISBNs
through the creation of actionable ISBNs, or ISBN-As [38], which may
serve as a model for actionable content hashes.

What it Means to Preserve Data

Our results indicate that reference rot poses an existential threat to
published biodiversity datasets. We’ve seen that the use of content-based
identifiers can effectively address the issue of reference rot. However,
identifiers are of little use in a vacuum. An identifier can only be useful for
data retrieval when combined with a resolver to associate identifiers with
locations and a database to retrieve the dataset at the associated location
([24]). Thus, we need to address how resolvers and databases might be
organized to accommodate content-based identifiers in order to fully realize
long-term data preservation. In this context, we define data preservation

as the continued capacity for datasets to be reliably referenced and
retrieved in their original form even as the global digital biodiversity
network evolves over time.

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

We have shown that software agents such as Preston can be used to
collect datasets and their provenance over time while maintaining
content-addressability; all that is needed to ensure proper data preservation
are a dedicated repository and a well-known, publicly available search
index to map content-based identifiers to datasets located in the repository.
In practice, repositories and search indexes (and potentially software
agents such as Preston deployments) can be co-located; examples include
Zenodo and the Internet Archive, although they impose some limitations
that may restrict file size, number of files, and the amount of information
that can be indexed ([41], [1]). These existing information systems may

serve as models for long-term biodiversity information systems. 613

These requirements help to ensure that biodiversity data remain FAIR 614
(Findable, Accessible, Interoperable, and Reusable) ([40]). Findability is 615
achieved through the publishing of provenance logs which thoroughly 616
describe what datasets are and where they originated from. The 617
amenability of the content-based identification paradigm to the operation 618
of independent distributed repositories strengthens accessibility by 619
preventing the failure of a single data repository from inhibiting future 620
data access (4). Content-based identification also allows for interoperability 621
due to the absence of any central authority to administrate data access; a 622
content hash computed from a dataset is guaranteed to match the hash 623
computed by any other agent using the same dataset. Finally, and 624
particularly relevant to this paper’s purpose, reusability is strengthened by 625
enhancing the retrievability of referenced datasets and allowing users to 626
verify that a retrieved dataset exactly matches that which was referenced. 627

Conclusions 628

Although reference rot is resulting in a steady decline in the reliability of 629
our digital biodiversity record, realistic solutions are available to address 630
the root causes of the issue. Content drift can be eliminated altogether by 631
changing the way we reference datasets, from using location-based 632
identifiers to ones that are content-based. Meanwhile, the online 633
biodiversity data networks can be made far more resilient to link rot if 634
distributed observation and archival techniques are used to capture 635
incremental changes to the data record so that references can remain valid 636
even when online datasets are updated, removed, or relocated. 637

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

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

Acknowledgments

We thank ...

References

1. I. Archive. Uploading - a basic guide, 2019. Accessed: 2019-12-04.
2. D. Bearman. Archival strategies. *The American Archivist*, 58(4):380–413, sep 1995.

3. T. Berners-Lee, R. T. Fielding, and L. M. Masinter. Uniform Resource Identifier (URI): Generic Syntax. RFC 3986, Jan. 2005.
4. T. Berners-Lee, L. M. Masinter, and M. P. McCahill. Uniform Resource Locators (URL). RFC 1738, Dec. 1994.
5. M. J. Costello, P. Bouchet, G. Boxshall, K. Fauchald, D. Gordon, B. W. Hoeksema, G. C. B. Poore, R. W. M. van Soest, S. Stöhr, T. C. Walter, B. Vanhoorne, W. Decock, and W. Appeltans. Global coordination and standardisation in marine biodiversity through the world register of marine species (WoRMS) and related databases. *PLoS ONE*, 8(1):e51629, jan 2013.
6. DataONE. Dataone citation guidelines, 2012. Accessed: 2019-12-04.
7. E. B. Davis and D. Schmidt. *Guide to Information Sources in the Botanical Sciences*. Vol. 2nd ed. Reference Sources in Science and Technology. Englewood, Colo: Libraries Unlimited., 1996.
8. J. L. Edwards. Interoperability of biodiversity databases: Biodiversity information on every desktop. *Science*, 289(5488):2312–2314, sep 2000.
9. C. S. Elton. *Animal ecology*. Macmillan Co., 1927.
10. G. T. G. B. I. Facility. What is gbif?, 2019. Accessed: 2019-12-04.
11. I. D. Foundation. Doi handbook, 2018. Accessed: 2019-12-04.
12. E. Garfield, I. H. Sher, and R. J. Torpie. *The Use of Citation Data in Writing the History of Science*. Institute for Scientific Information Inc Philadelphia PA, 1964.

13. GBIF.org. Gbif citation guidelines, 2019. Accessed: 2019-12-04.
14. A. Güntsch, R. Hyam, G. Hagedorn, S. Chagnoux, D. Röpert, A. Casino, G. Droege, F. Glöckler, K. Gödderz, Q. Groom, J. Hoffmann, A. Holleman, M. Kempa, H. Koivula, K. Marhold, N. Nicolson, V. S. Smith, and D. Triebel. Actionable, long-term stable and semantic web compatible identifiers for access to biological collection objects. *Database*, 2017, jan 2017.
15. J. Hortal, F. de Bello, J. A. F. Diniz-Filho, T. M. Lewinsohn, J. M. Lobo, and R. J. Ladle. Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual Review of Ecology, Evolution, and Systematics*, 46(1):523–549, dec 2015.
16. iDigBio.org. idigbio citation guidelines, 2016. Accessed: 2019-12-04.
17. M. Klein, H. V. de Sompel, R. Sanderson, H. Shankar, L. Balakireva, K. Zhou, and R. Tobin. Scholarly context not found: One in five articles suffers from reference rot. *PLoS ONE*, 9(12):e115253, dec 2014.
18. T. Kuhn and M. Dumontier. Making digital artifacts on the web verifiable and reliable. *IEEE Transactions on Knowledge and Data Engineering*, 27(9):2390–2400, sep 2015.
19. A. Matsunaga, R. Figueiredo, A. Thompson, G. Traub, R. Beaman, and J. A. Fortes. Integrated digitized biocollections (idigbio) cyberinfrastructure status and futures. In *TDWG 2013 ANNUAL CONFERENCE*, 2013.

20. W. Michener, D. Vieglais, T. Vision, J. Kunze, P. Cruse, and G. Janée. DataONE: Data observation network for earth: Preserving data and enabling innovation in the biological and environmental sciences. *D-Lib Magazine*, 17(1/2), jan 2011.
21. P. Mockapetris. Domain names - concepts and facilities. RFC 1034, Nov. 1987.
22. NIST. Descriptions of sha-256, sha-384, and sha-512, 2001. Accessed: 2019-12-04.
23. L. M. Page, B. J. MacFadden, J. A. Fortes, P. S. Soltis, and G. Riccardi. Digitization of biodiversity collections reveals biggest data on biodiversity. *BioScience*, 65(9):841–842, aug 2015.
24. N. Paskin. Toward unique identifiers. *Proceedings of the IEEE*, 87(7):1208–1227, jul 1999.
25. N. Paskin. Digital object identifier (DOI®) system. In *Encyclopedia of Library and Information Sciences, Third Edition*, pages 1586–1592. CRC Press, dec 2009.
26. J. Poelen, M. Elliott, I. Alzuru, and P. Patel. Preston: a biodiversity dataset tracker, Sept. 2018.
27. J. H. Poelen. A biodiversity dataset graph: Biodiversity Heritage Library (BHL), June 2019.
28. J. H. Poelen. A biodiversity dataset graph: DataONE, July 2019.
29. J. H. Poelen. A biodiversity dataset graph: GBIF, iDigBio, BioCAsE, Oct. 2019.

30. J. H. Poelen. Biodiversity Dataset Archive, Oct. 2019.
31. J. H. Poelen. A biodiversity dataset graph: Bhl, Oct. 2019.
32. J. H. Poelen. A biodiversity dataset graph: Dataone, Oct. 2019.
33. J. H. Poelen. To connect is to preserve: on frugal data integration and preservation solutions, Jun 2019.
34. C. Rinaldo and C. Norton. BHL, the biodiversity heritage library: An expanding international collaboration. *Nature Precedings*, aug 2009.
35. T. Robertson, M. Döring, R. Guralnick, D. Bloom, J. Wiczorek, K. Braak, J. Otegui, L. Russell, and P. Desmet. The GBIF integrated publishing toolkit: Facilitating the efficient publishing of biodiversity data on the internet. *PLoS ONE*, 9(8):e102623, aug 2014.
36. B. Trask. Principles of content addressing.
<https://bentrask.com/?q=hash://sha256/98493caa8b37eaa26343bbf73f232597a3ccda20498563327a4c3713821df892>,
2015. Accessed: 2019-12-04.
37. T. J. Vision. Open data and the social contract of scientific publishing. *BioScience*, 60(5):330–331, may 2010.
38. A. Weissberg. The identification of digital book content. *Publishing Research Quarterly*, 24(4):255–260, nov 2008.
39. J. Wiczorek, D. Bloom, R. Guralnick, S. Blum, M. Döring, R. Giovanni, T. Robertson, and D. Vieglais. Darwin core: An

evolving community-developed biodiversity data standard. *PLoS ONE*, 7(1):e29715, jan 2012.

40. M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. T. Evelo, R. Finkers, A. Gonzalez-Beltran, A. J. Gray, P. Groth, C. Goble, J. S. Grethe, J. Heringa, P. A. 't Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S. J. Lusher, M. E. Martone, A. Mons, A. L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.-A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M. A. Swertz, M. Thompson, J. van der Lei, E. van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, and B. Mons. The FAIR guiding principles for scientific data management and stewardship. *Scientific Data*, 3(1), mar 2016.
41. Zenodo. General policies, 2019. Accessed: 2019-12-04.

Figures

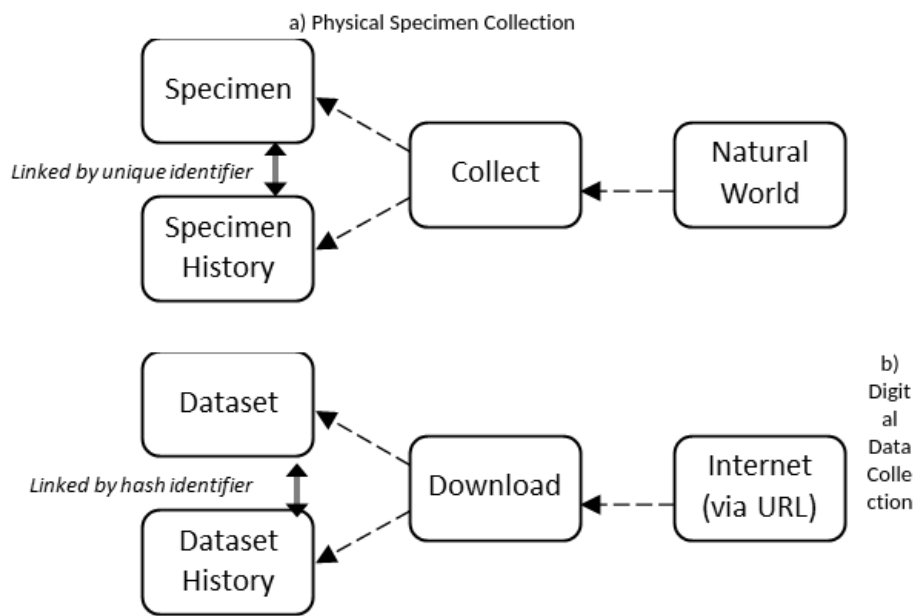


Figure 1. Reliable record keeping for digital datasets (b) can be achieved in an analogous way to current practices in record keeping for physical specimens (a). Biologists collect physical specimens from the natural world, thoroughly document the process, then store the specimens in facilities equipped 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.

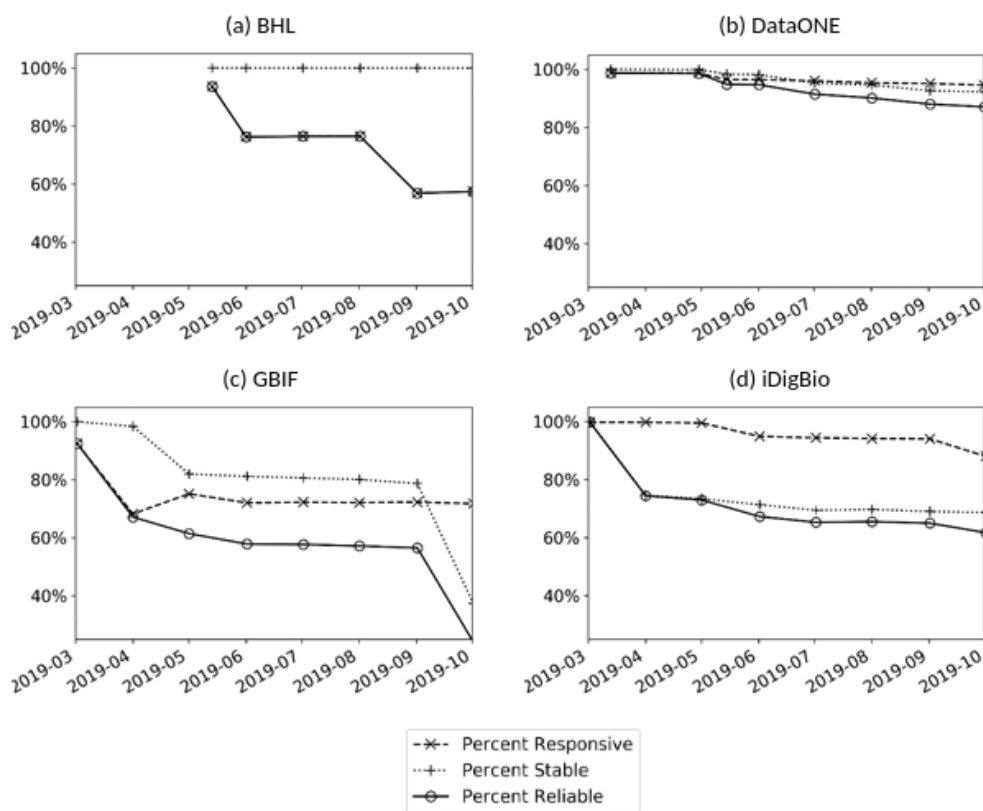


Figure 2. Overall responsiveness, stability, and reliability from March 2019 to October 2019 as a percentage of URLs that exhibit each indicator in a) BHL, b) DataONE, c) GBIF, and d) iDigBio.

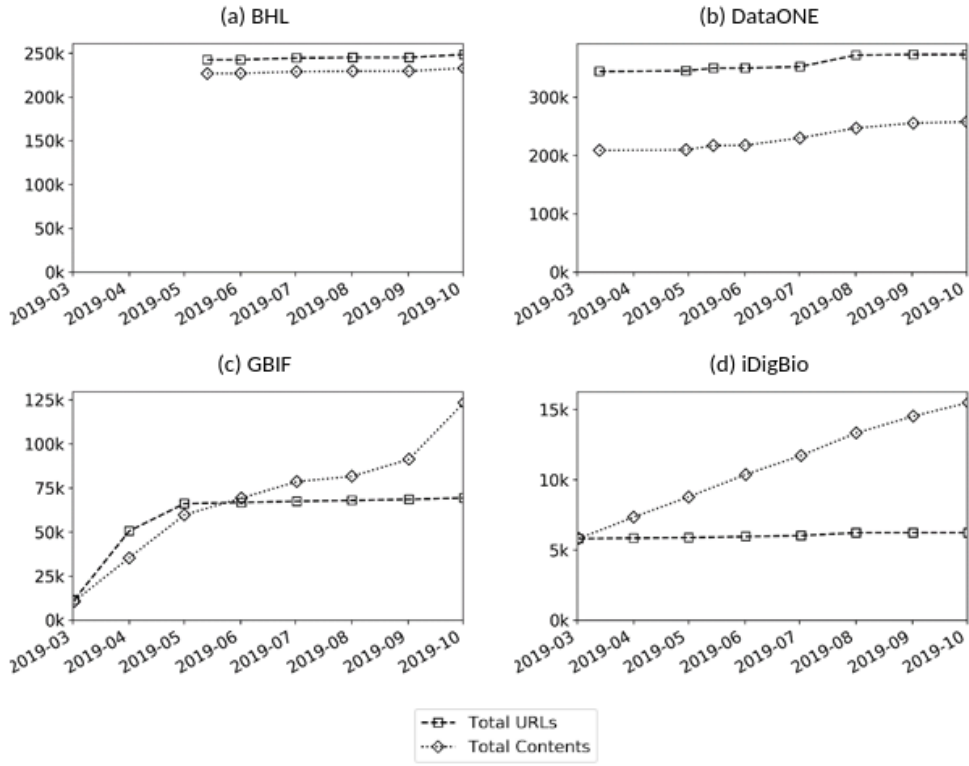


Figure 3. Total number of URLs and unique contents observed from March 2019 to October 2019 for a) BHL, b) DataONE, c) GBIF, and d) iDigBio.

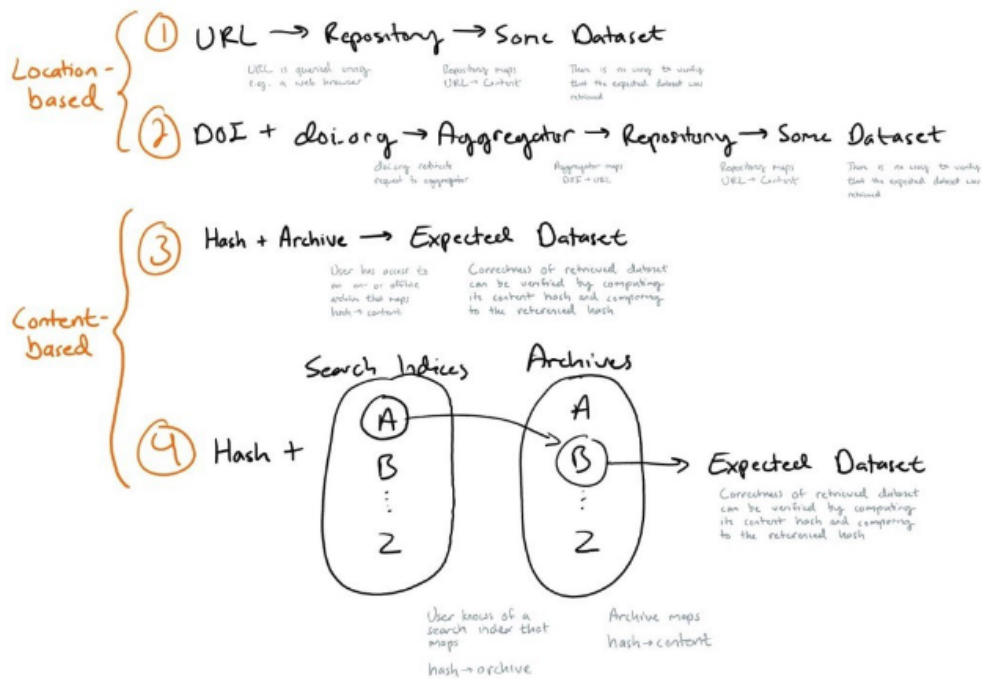


Figure 4. Visualization of content resolution for location- versus content-based identifiers. 1) URLs point to a known location of a dataset, but do not guarantee either the presence or authenticity of the retrieved dataset; 2) the use of a DOI that resolves to a URL adds a layer of redirection; 3) A 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.

Hash Archive (beta)

URL or hash:

Sources for [hash://sha256/b83cf099449dae3f633af618b19d05013953e7a1d7d97bc5ac01afd7bd9abe5d](https://archive.org/download/biodiversity-dataset-archives/data.zip/data/b8/3c/b83cf099449dae3f633af618b19d05013953e7a1d7d97bc5ac01afd7bd9abe5d/)

- [Search for this hash on Google](#)
- [Search for this hash on DuckDuckGo](#)
- [Search for this block on IPFS](#)
- [Check this hash on VirusTotal](#)
- [Other useful sources...?](#)

Active as of November 5th, 2019

<https://archive.org/download/biodiversity-dataset-archives/data.zip/data/b8/3c/b83cf099449dae3f633af618b19d05013953e7a1d7d97bc5ac01afd7bd9abe5d/>^[1]

Active as of October 8th, 2019

<https://deeplinker.bio/b83cf099449dae3f633af618b19d05013953e7a1d7d97bc5ac01afd7bd9abe5d/>^[1]

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.