#### Conference

4th Annual Digital Data in Biodiversity Research Conference June 1-3 2020

#### Title

Reliable dataset identifiers are essential building blocks for reproducible research

#### Type

Oral presentation, Discussion

#### **Authors**

Jorrit Poelen, Michael Elliott, Jose Fortes

#### **Abstract**

Despite increased use of digital biodiversity data in research, reliable methods to identify datasets are not widely adopted. While commonly used location-based dataset identifiers such as URLs help to easily download data today, additional identification schemes are needed to ensure long term access to datasets.

We propose to augment existing location- and DOI-based identification schemes with cryptographic content-based identifiers. These content-based identifiers can be calculated from the datasets themselves using available cryptographic hashing algorithms (e.g., sha256). These algorithms take only the digital content as input to generate a unique identifier without needing a centralized identification administration. The use of content-based identifiers is not new, but a re-application of change management techniques used in the popular version control system "git".

We show how content-based identifiers can be used to version datasets, to track the dataset locations, to monitor their reliability, and to efficiently detect dataset changes. We discuss the results of using our approach on datasets registered in GBIF and iDigBio from Sept 2018 to May 2020. Also, we propose how reliable, decentralized, dataset indexing and archiving systems can be devised. Lastly, we outline a modification to existing data citation practices to help work towards more reproducible and reusable research workflows.

# Reliable dataset identifiers are essential building blocks for reproducible research

1 June 2020, 4th Annual Digital Data in Biodiversity Research, see <u>"Reliable dataset identifiers are essential building blocks for reproducible research"</u>.

Jorrit Poelen, Ronin Institute and independent software engineer;

Michael Elliott, Advanced Computing and Information Systems Laboratory (ACIS) Department of Electrical and Computer Engineering, University of Florida;

Jose Fortes, Advanced Computing and Information Systems Laboratory (ACIS) Department of Electrical and Computer Engineering, University of Florida, Gainesville, FL

This work is funded in part by grant <u>NSF OAC 1839201</u> "EAGER: Towards the Web of Biodiversity Knowledge: Understanding Data Connectedness to Improve Identifier Practices" from the **National Science Foundation**.

## Reliable Dataset References

Michael Elliott, Jorrit H. Poelen, José A. B. Fortes





What does an unreliable reference look like?





■ An example reference using a URL

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.





An example reference using a URL

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.





An example reference using a URL

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





An example reference using a URL

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





Get data

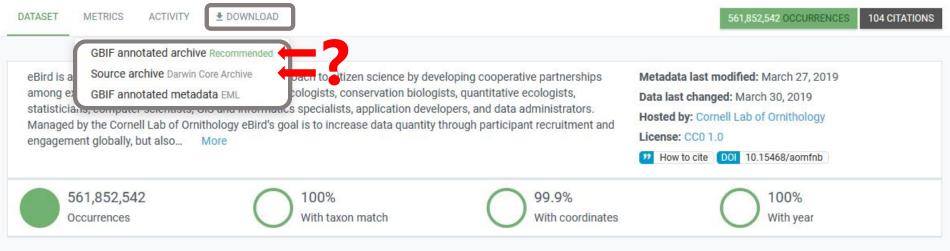
About

OCCURRENCE DATASET | REGISTERED SEPTEMBER 16, 2010

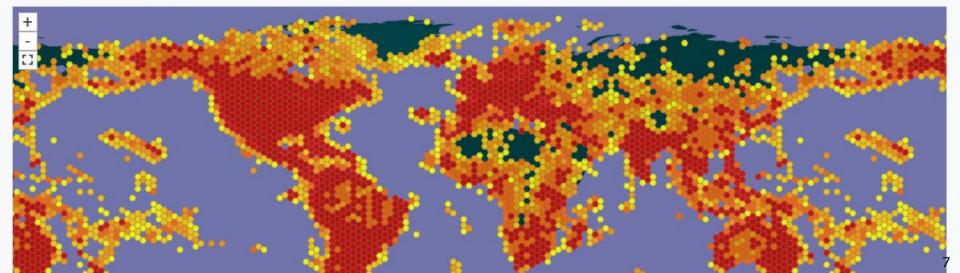
#### EOD - eBird Observation Dataset

Published by Cornell Lab of Ornithology

Tim Levatich • Francisco Padilla • 🖂 Jeff Gerbracht



#### 561,766,080 GEOREFERENCED RECORDS



Get data

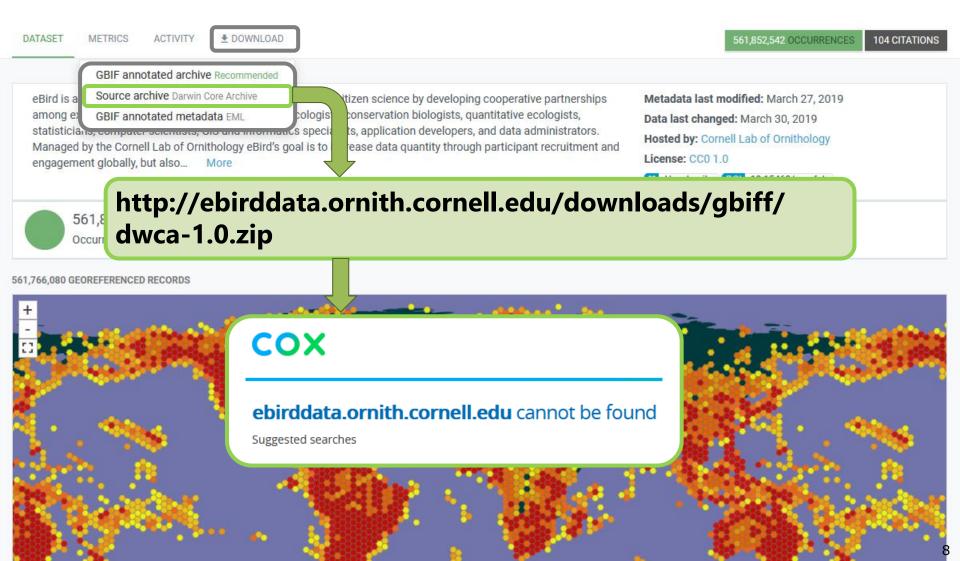
About

OCCURRENCE DATASET | REGISTERED SEPTEMBER 16, 2010

#### EOD - eBird Observation Dataset

Published by Cornell Lab of Ornithology

Tim Levatich • Francisco Padilla • 🖂 Jeff Gerbracht



Get data

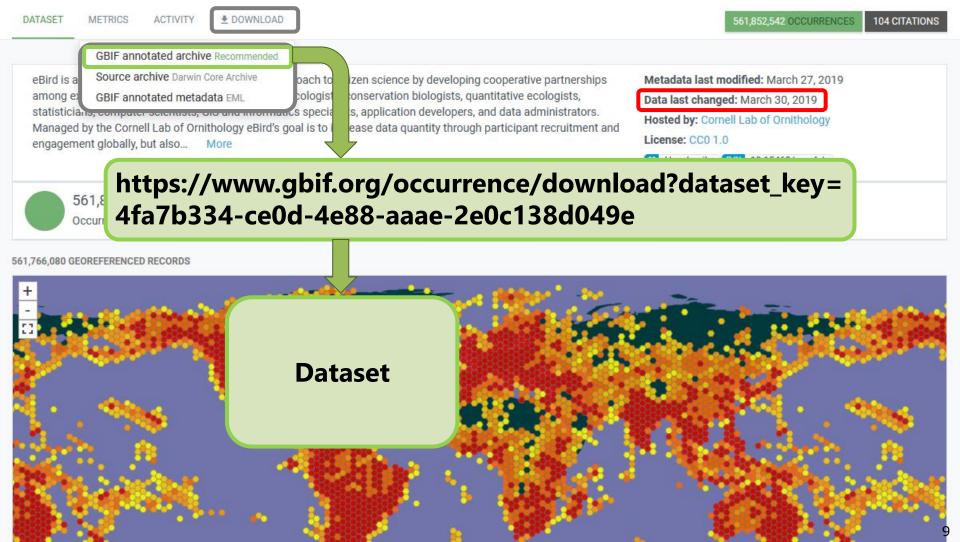
About

OCCURRENCE DATASET | REGISTERED SEPTEMBER 16, 2010

#### EOD - eBird Observation Dataset

Published by Cornell Lab of Ornithology

Tim Levatich • Francisco Padilla • 🖂 Jeff Gerbracht



About

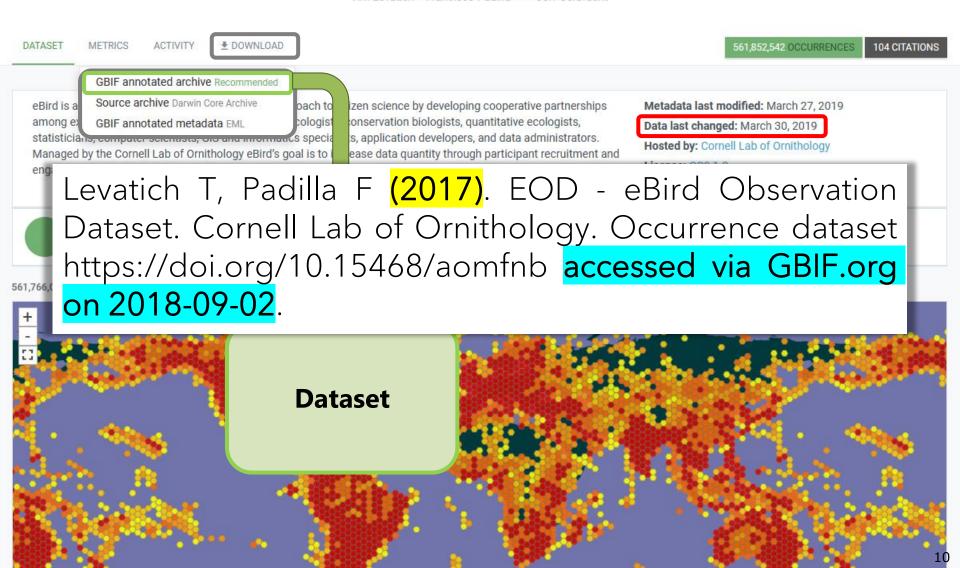
Get data

OCCURRENCE DATASET | REGISTERED SEPTEMBER 16, 2010

#### EOD - eBird Observation Dataset

Published by Cornell Lab of Ornithology

Tim Levatich • Francisco Padilla • 🖂 Jeff Gerbracht



A reference is **reliable** if it both

- Allows continued access to what was referenced
- Only identifies what was referenced

Location-based identifiers, such as URLs, tend to be unreliable

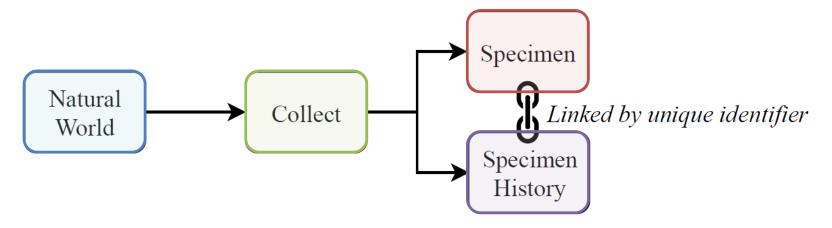




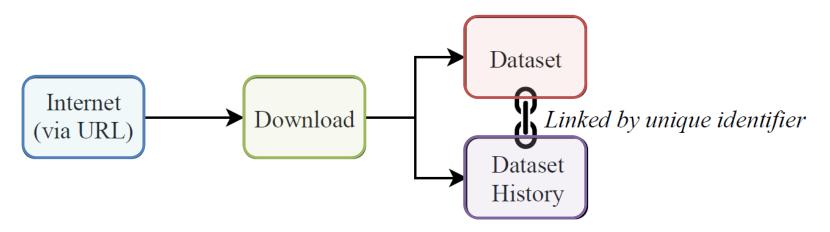
How can we reliably reference datasets served at unreliable URLs?







(a) Physical specimen collection



(b) Digital data collection





- Content-based identifiers can be used to identify datasets
- Cryptographic hash functions can produce unique content-based identifiers for digital datasets

Content-based identifier for the 2017 eBird dataset:

hash://sha256/<mark>29d30b566f924355a383b13cd48c3aa239</mark> d4\cba0a\5f4ccfc2930289b88b43c

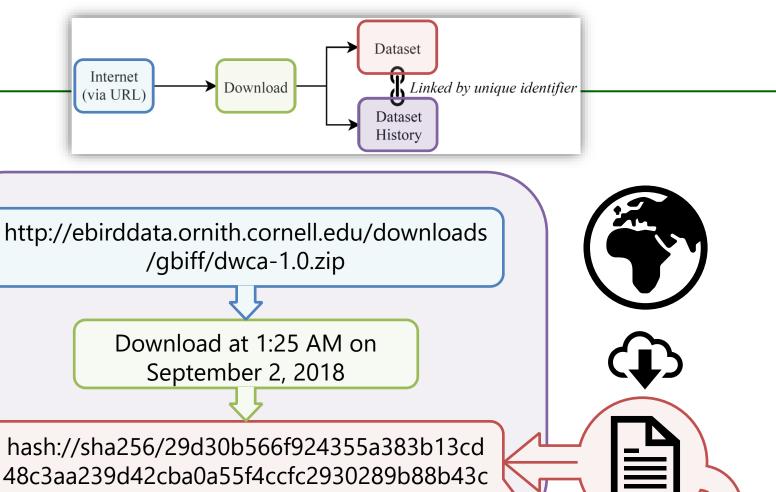


Content-based identifier for the 2019 eBird dataset:

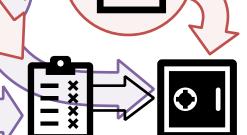
hash://sha256/ec3ff57cb48d5c41b77b5d1075738b40f5 98a900e8be56e7645e5a24013dffc4







hash://sha256/29d30b566f924355a383b13cd



hash://sha256/c253a5311a20c2fc082bf9bac8 7a1ec5eb6e4e51ff936e7be20c29c8e77dee55





## Our example eBird reference:

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.

## Add a content-based identifier

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





## Hash Archive (beta)

URL or hash:

hash://sha256/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c

Lookup

#### Sources for hash://sha256/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c

- · 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 May 27th, 2020

https://zenodo.org/record/3858251/files/dwca-1.0.zip?download=1[^]

Active as of May 27th, 2020

https://archive.org/download/biodiversity-dataset-archives/data.zip/data/29/d3/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c [^]

Active as of May 26th, 2020

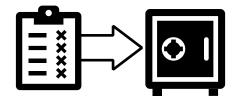
https://deeplinker.bio/29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930289b88b43c[^]

Active as of February 12<sup>nd</sup>, 2020

https://archive.org/download/biodiversity-dataset-archives/data.zip/data%2F29%2Fd3%2F29d30b566f924355a383b13cd48c3aa239d42cba0a55f4ccfc2930 289b88b43c[^]

Obsolete; last seen March 9th, 2019

http://ebirddata.ornith.cornell.edu/downloads/gbiff/dwca-1.0.zip[^]





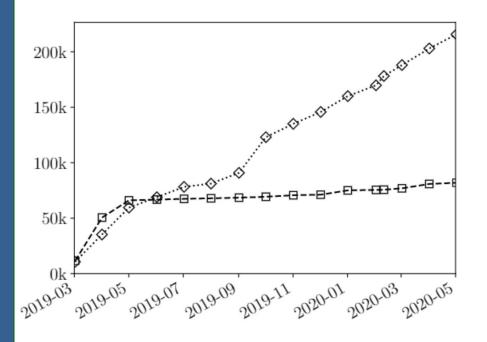


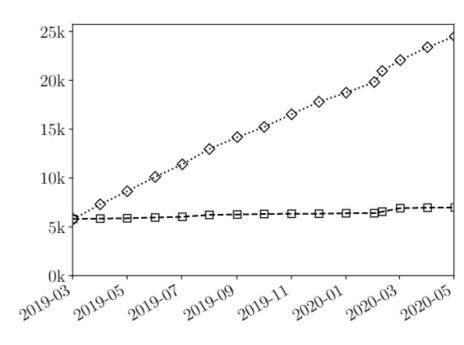
- In order to have reliable references,
  - 1) Datasets must be addressable and retrievable using content-based identifiers (rather than, e.g., URLs)
  - 2) Datasets that could be used in the future should be archived
  - Openly accessible registries must exist to lookup locations for datasets referenced by content-based identifiers
  - 4) Agents must exist to collect datasets, record their provenance (i.e., history record), deposit them into archives, and share their new locations with public registries





- Why are reliable references important?
  - The reliability of URLs as references decreases over time





(c) GBIF

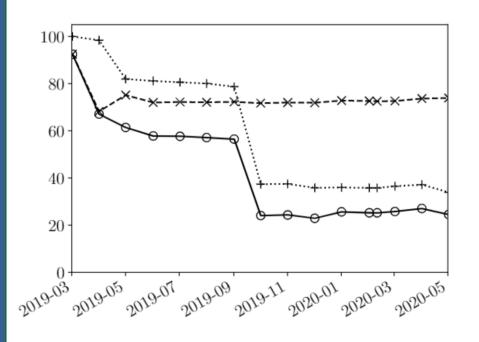
Total URLs Total Contents (Datasets)

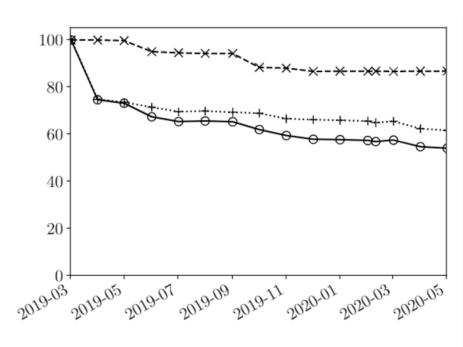
(d) iDigBio



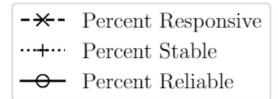


- Why are reliable references important?
  - The reliability of URLs as references decreases over time





(c) GBIF



(d) iDigBio





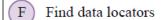
# Thanks for Listening!

Based on our paper, "Toward Reliable Biodiversity Dataset References" (Preprint: <a href="https://ecoevorxiv.org/mysfp">https://ecoevorxiv.org/mysfp</a>)

The research was funded in part by a grant (NSF OAC 1839201) from the National Science Foundation and the AT&T Foundation.





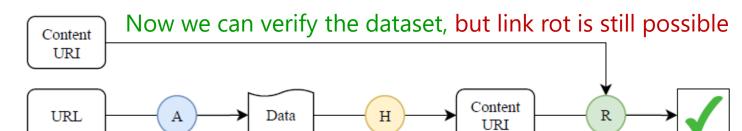


- A Access located data
- H Compute content hash
- R Reproducibility check
- ? Cannot verify content
- ✓ Can verify content

## Can't verify whether the retrieved dataset is what was referenced

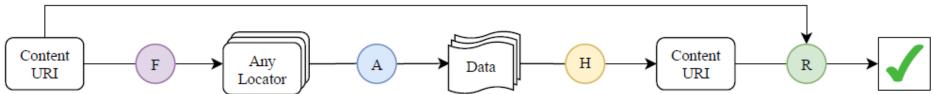


(a) URL reference



(b) URL reference with content hash

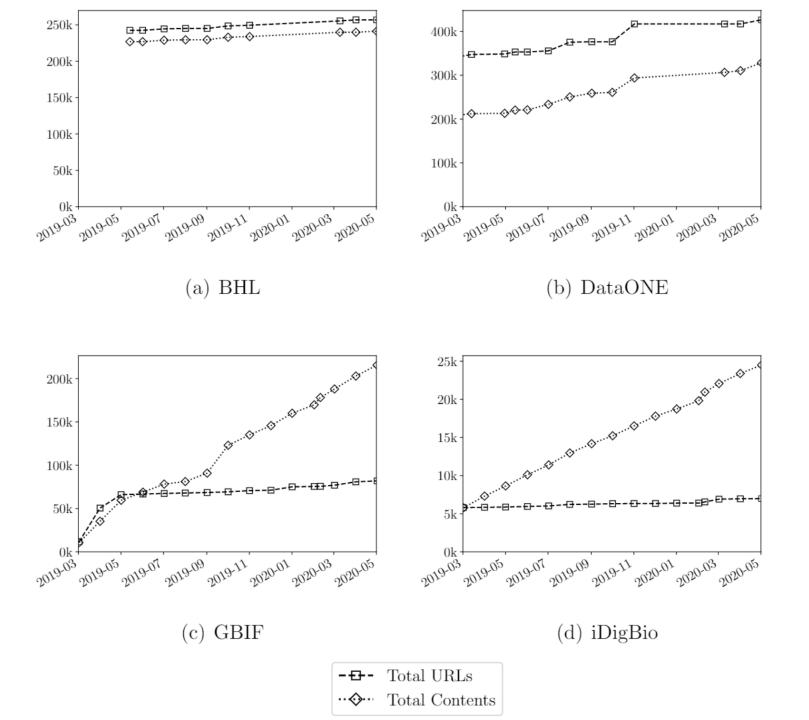
Now we can reliably find and verify the dataset

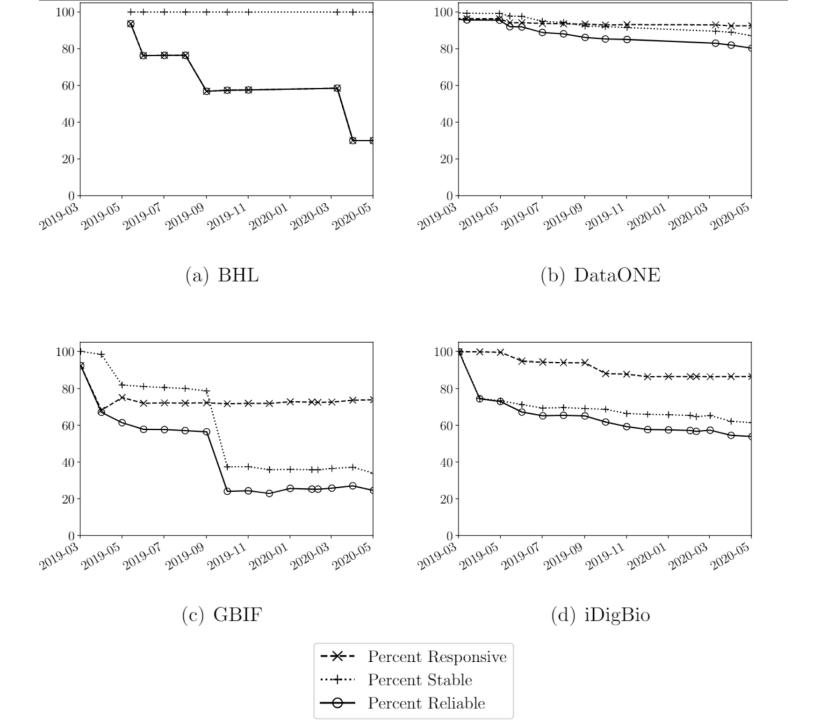


(c) Content URI reference









## Agenda (see <a href="https://tinyurl.com/reliable-data">https://tinyurl.com/reliable-data</a>)

- 1. Primer Questions (1 min)
- 2. Michael presents "Reliable Dataset References" (10 mins)
- 3. Revisit Primer Questions (1 min)
- 4. Discussion (45 12 = 33 mins)
- 5. Sign-up for <u>post-discussion hallway conversation</u> at <a href="https://ufl.zoom.us/j/98346741161">https://ufl.zoom.us/j/98346741161</a> at 4PM EST (BYO coffee and cookies!)

## **Primer Questions**

How do you make sure to keep your data around?

How do you (reliably) reference data in publications?

How do **we** identify **our** digital data so that it can be kept for the next 50 years?

Michael presents "Reliable Dataset References"

#### **Revisit Primer Questions**

How do **you** make sure to keep **your** data around?

How do **you** (reliably) reference data in publications?

How do **we** identify **our** digital data so that it can be kept for the next 50 years?

## Post-Discussion Hallway Conversation Signup:

Please add your name / affiliation / content info via <a href="https://tinyurl.com/reliable-data">https://tinyurl.com/reliable-data</a> and join us at <a href="https://ufl.zoom.us/j/98346741161">https://ufl.zoom.us/j/98346741161</a> at 4PM EST for a casual hallway conversation.

(edited: removed email addresses to reduce exposure)

Name	Affiliation	Contact Info
Jorrit Poelen	Ronin Institute	xxx
Michael Elliott	University of Florida	xxx
Beth Chambers	William & Mary	xxx
Gil Nelson	iDigBio	xxx
Jeff Gerbracht	Cornell Lab of Ornithology	xxx

Jose Fortes	UF ACIS/iDigBio	xxx
Randy Singer	University of Michigan	xxx
Gabriel	Kamener	xxx
Ali Krzton	Auburn University	xxx
Deborah Paul	Florida State University	xxx

#### References

Levatich T, Padilla F (2019). EOD - eBird Observation Dataset. Cornell Lab of Ornithology. Occurrence dataset <a href="https://doi.org/10.15468/aomfnb">https://doi.org/10.15468/aomfnb</a> accessed via GBIF.org on 2019-04-08 <a href="https://sha256/ec3ff57cb48d5c41b77b5d1075738b40f598a900e8be56e7645e5a24013dffc4">https://sha256/ec3ff57cb48d5c41b77b5d1075738b40f598a900e8be56e7645e5a24013dffc4</a>

Elliott, M. J., Poelen, J. H., & Fortes, J. (2020, January 3). Toward Reliable Biodiversity Dataset References. <a href="https://doi.org/10.32942/osf.io/mysfp">https://doi.org/10.32942/osf.io/mysfp</a>. For related GBIF-mediated forum discussion: <a href="https://discourse.gbif.org/t/toward-reliable-biodiversity-dataset-references/1637/3">https://discourse.gbif.org/t/toward-reliable-biodiversity-dataset-references/1637/3</a>.

Trask B. 2015. Principles of content addressing.

https://bentrask.com/?q=hash://sha256/98493caa8b37eaa26343bbf73f232597a3ccda20498563327a4c3713821df892. Accessed: 2019-12-04

NSF OAC 1839201 "EAGER: Towards the Web of Biodiversity Knowledge: Understanding Data Connectedness to Improve Identifier Practices" from the **National Science Foundation**.