# MAKING SMITHSONIAN OPEN ACCESS ACCESSIBLE WITH PYTHON AND DASK

## Mike Trizna

Smithsonian OCIO Data Science Lab

May 4, 2021 | *csv,conf,v6*

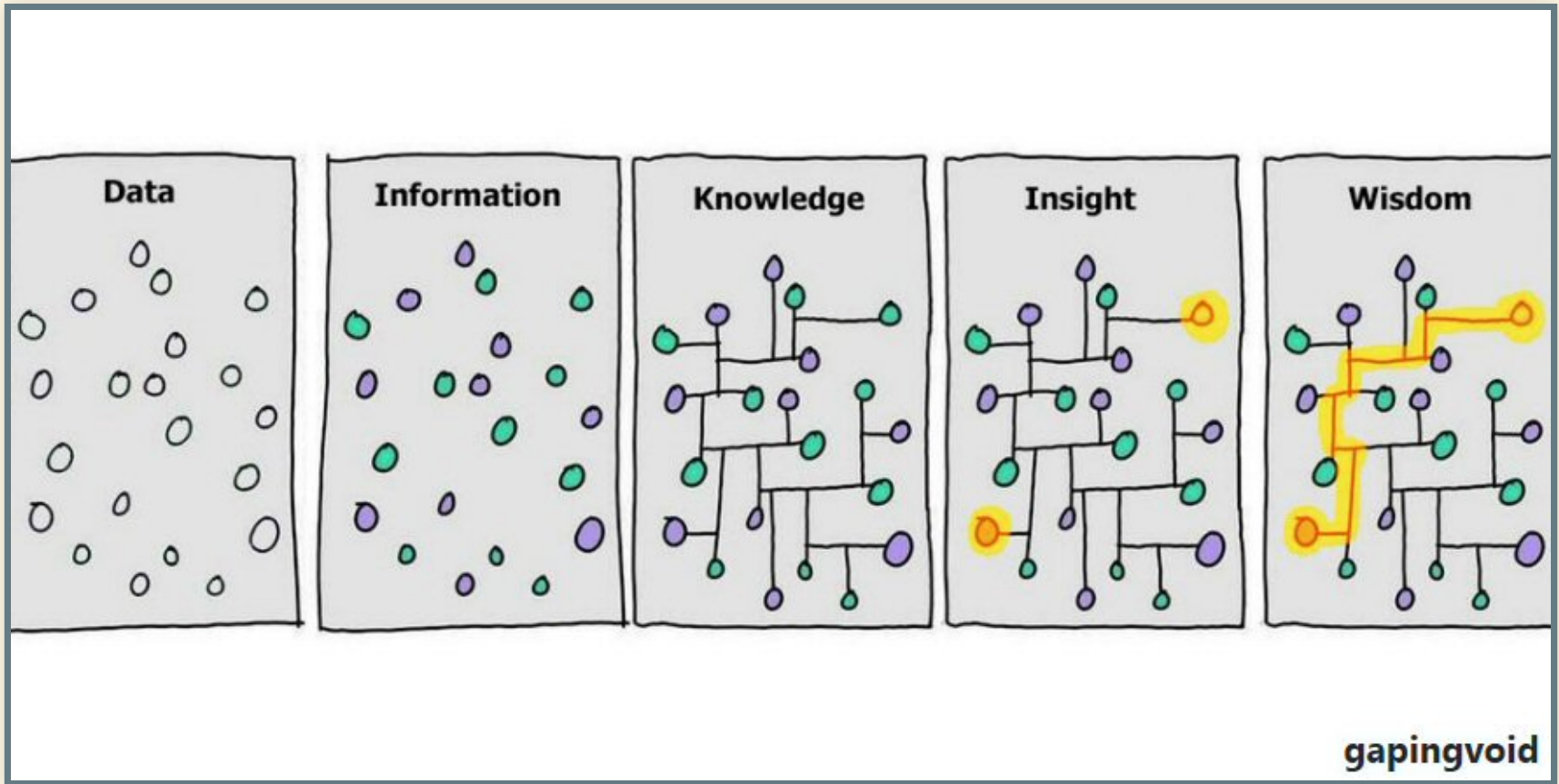# WHAT IS THE SMITHSONIAN INSTITUTION?

Yes, there are the museums (19 of them, mostly in Washington, DC), but we also have 21 libraries and archives, 9 research centers ... and a zoo.

# SMITHSONIAN MISSION

Founded in 1846 from the bequest of Englishman James Smithson with the condition:

> *"under the name of the Smithsonian Institution, an establishment for the* ***increase and diffusion of knowledge.***"
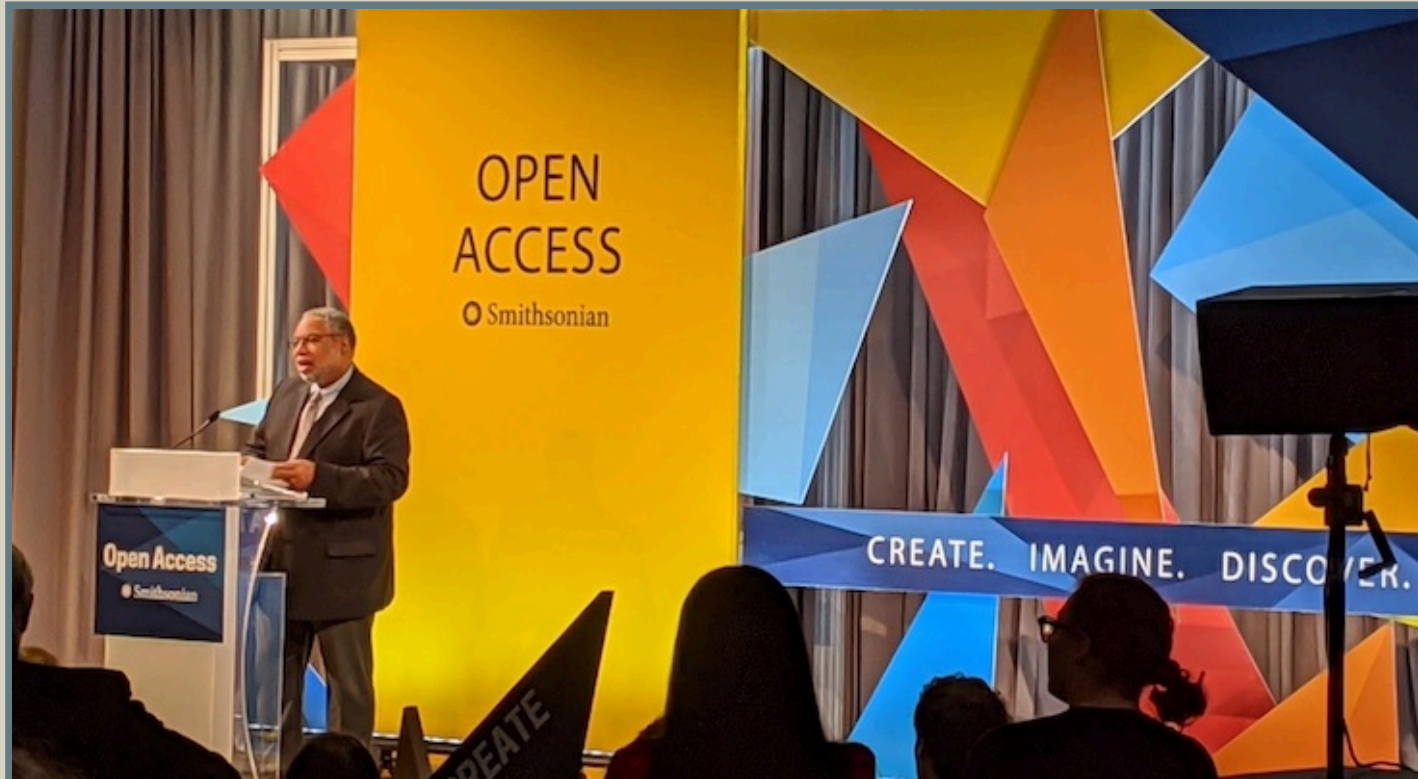
The Smithsonian has been increasing and diffusing *Knowledge* since 1846, but what about all of that *Data*?

All of that *data* and *info* that fed into *knowledge, insight, and wisdom* were dutifully cataloged and stored.

# LAUNCH EVENT

## February 25, 2020

# SI OPEN ACCESS RELEASE

Of the Smithsonian's 155 million objects, 2.1 million library volumes and 156,000 cubic feet of archival collections:

- 2.8 million 2-D and 3-D images
- Over 17 million collection metadata objects

# TERMS OF USE

Before February 2020, all Smithsonian museums and units made their data searchable and sometimes able to download, but through each individual unit. Many different use agreements.

SI Open Access put all media and metadata in one place, and all Open Access media is CC0.

# SO HOW CAN ALL OF THIS DATA BE ACCESSED?

I will cover 3 different ways.

All 3 share metadata records in the same deeply-nested JSON structure.

# WEB API

🔗: http://edan.si.edu/openaccess/apidocs/

- API Key needed (but free and painless to register)
- Great for getting a feel for record structure

# HOWEVER YOU WILL QUICKLY RUN INTO LIMITATIONS

- Records are extensively indexed, but can only search indexed fields.
- Row limit of 1000 per API call

# EXAMPLE: HERBARIUM SHEET MERCURY DETECTOR



**Biodiversity Data Journal 5: e21139**
doi: 10.3897/BDJ.5.e21139

**Research Article**

**Applications of deep convolutional neural networks to digitized natural history collections**

Eric Schuettpelz‡, Paul B. Frandsen§, Rebecca B. Dikow§, Abel Brown|, Sylvia Orli‡, Melinda Peters‡, Adam Metallo§, Vicki A. Funk‡, Laurence J. Dorr‡

‡ National Museum of Natural History, Smithsonian Institution, Washington, DC, United States of America
§ Office of the Chief Information Officer, Smithsonian Institution, Washington, DC, United States of America
| NVIDIA, Santa Clara, CA, United States of America



2017 paper that described building a machine learning model to detect herbarium sheets that had been stained with mercury.

https://doi.org/10.3897/bdj.5.e21139

# EXAMPLE: HERBARIUM SHEET MERCURY DETECTOR

I wanted to create a new model on same dataset (2017 is ancient history in Machine Learning)

All training images are shared on Figshare, but photos are resized and I wanted original metadata

# EXAMPLE: HERBARIUM SHEET MERCURY DETECTOR

Unfortunately the "barcode" term from the supplementary materials is not an indexed field.

# FULL DATASET SOURCES: GITHUB AND AWS

- **AWS S3** for all metadata and images
- **GitHub** for versioned metadata

# HOW IT'S PACKAGED

- Files are serialized as line-delimited JSON and compressed with bzip2.
- Directories are organized by owning unit and files are distributed by first two characters of content serialization hash.

# HOW IT'S PACKAGED

# BUMMER, WE NEED TO GO THROUGH EVERY SINGLE FILE ONE AT A TIME?

If I'm looking across all units, that's 9,728 files to process!

# BUT THERE'S ACTUALLY A BENEFIT TO SO MANY FILES TO PROCESS

# WE CAN MULTITASK!

# ENTER DASK



https://dask.org/

# WHAT IS DASK?

Dask lets you set up a mini cluster on your machine ... or on an actual compute cluster

# DASK DASHBOARD

# WHAT IS DASK?

Dask is more well-known for parallel processing of DataFrames, but it also contains a really useful catch-all "Bag" type.

```python
import dask.bag as db
import json
s3_bag = db.read_text('s3://smithsonian-open-access/metadata/edan/*/*.txt',
                      storage_options={'anon': True}).map(json.loads)
gh_bag = db.read_text('~/Documents/OpenAccess/metadata/objects/*/*.txt.bz2',
                      compression='bz2').map(json.loads)
```

# EXAMPLE: HERBARIUM SHEET MERCURY DETECTOR

```
just_ids = (gh_bag.map(extract_ids)
                  .compute())
just_ids_df = pd.DataFrame(just_ids)
just_ids.head()
```

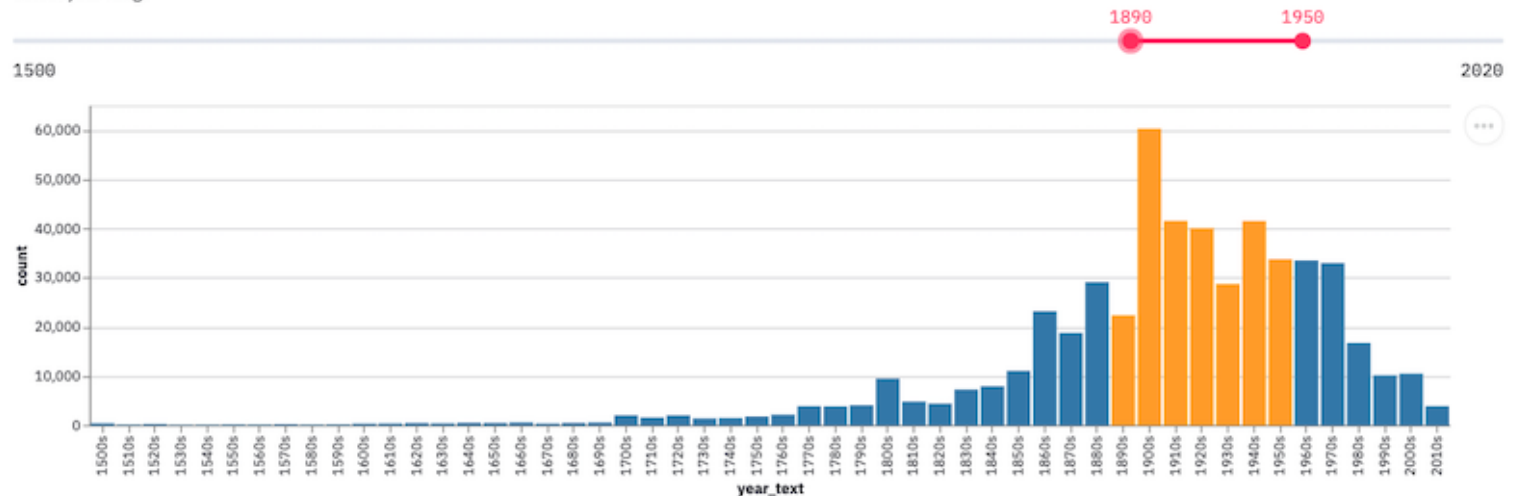| | edan_id | title | Barcode | specimen_guid | media_count | media_guid | ids_id | USNM Number |
|---|---|---|---|---|---|---|---|---|
| 0 | edanmdm-nmnhbotany_2095387 | Ageratum elachycarpum B.L. Rob. | 00512770 | http://n2t.net/ark:/65665/3d361b9b9-91d0-4a3d-... | 1.0 | http://n2t.net/ark:/65665/m3890da9ed-3b6d-40da... | NMNH-00512770 | 1404283 |
| 1 | edanmdm-nmnhbotany_2135634 | Horsfieldia bartlettii Merr. | 00513412 | http://n2t.net/ark:/65665/319872980-bd6d-409a-... | 1.0 | http://n2t.net/ark:/65665/m374b20557-3460-4d9d... | NMNH-00513412-000001 | 2275439 |
| 2 | edanmdm-nmnhbotany_2102106 | Hypericum crenulatum var. major Boiss. | 00588520 | http://n2t.net/ark:/65665/353447a40-80df-4892-... | 1.0 | http://n2t.net/ark:/65665/m3566c5f4e-c485-4c1b... | NMNH-00588520 | 129657 |
| 3 | edanmdm-nmnhbotany_2167183 | Asclepias brachystephana Engelm. ex Torr. in E... | 00588654 | http://n2t.net/ark:/65665/3cd0feeeb-abcd-415c-... | 1.0 | http://n2t.net/ark:/65665/m3330aaa7b-0613-4374... | NMNH-00588654-000001 | 18691 |
| 4 | edanmdm-nmnhbotany_2683657 | Waltheria indica L. | 00595768 | http://n2t.net/ark:/65665/3cbca3c96-e7ec-4eff-... | 2.0 | http://n2t.net/ark:/65665/m328e137c2-1b2e-4768... | NMNH-00595768 | 13147 |

# DIGITAL HUMANITIES INTERNSHIP PROJECT

## PATRICK MCMANUS FROM GEORGE MASON UNIVERSITY

# EXAMPLE: SEMANTIC CLUSTERING AMERICAN ART PAINTINGS



Full interactive notebook (through Binder) available at
https://github.com/sidatasciencelab/siopenaccess.

# QUESTIONS?

https://www.si.edu/openaccess

https://github.com/MikeTrizna/CSVConf2021_siopenacc