National Synchrotron Light Source II

**Brookhaven™**
National Laboratory

**U.S. DEPARTMENT OF**
**ENERGY**

# *Tiled* in the Context of Data and Metadata Services

Daniel Allan

Data Engineering Group Lead
Data Science and Systems Integration Program, NSLS-II
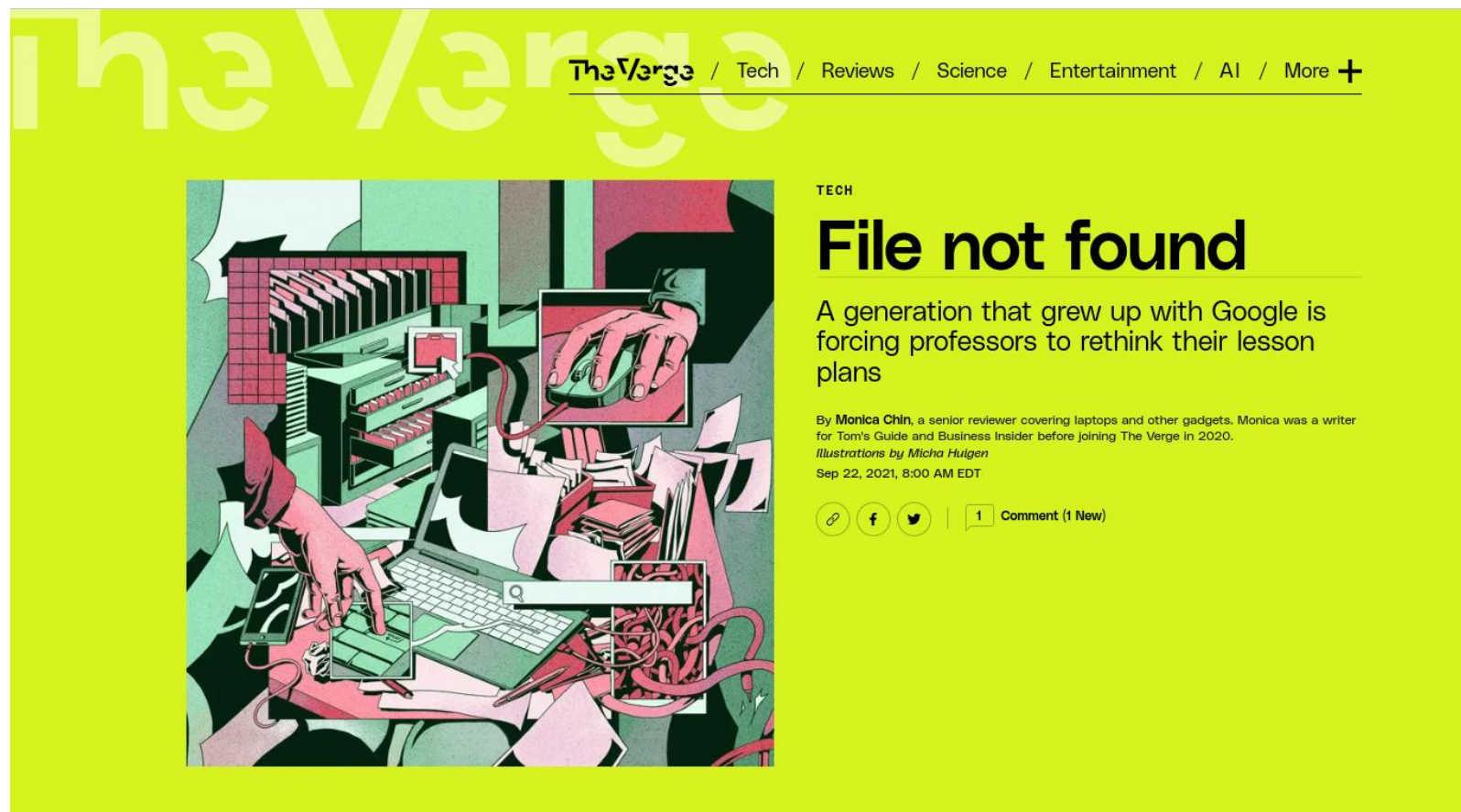
# There are many Data Services

DataFed

ArrayLake

Xpublish

h5grove

psana

Globus

HSDS

Tiled

datasette

SciCat

DVID

# Who Ordered a "Data Service"?

- Remove friction from data analysis at the small scale

- Enable new science at the large scale

- Make it easy for busy researchers to be F.A.I.R. (Findable, Accessible, Interoperable, Reusable)
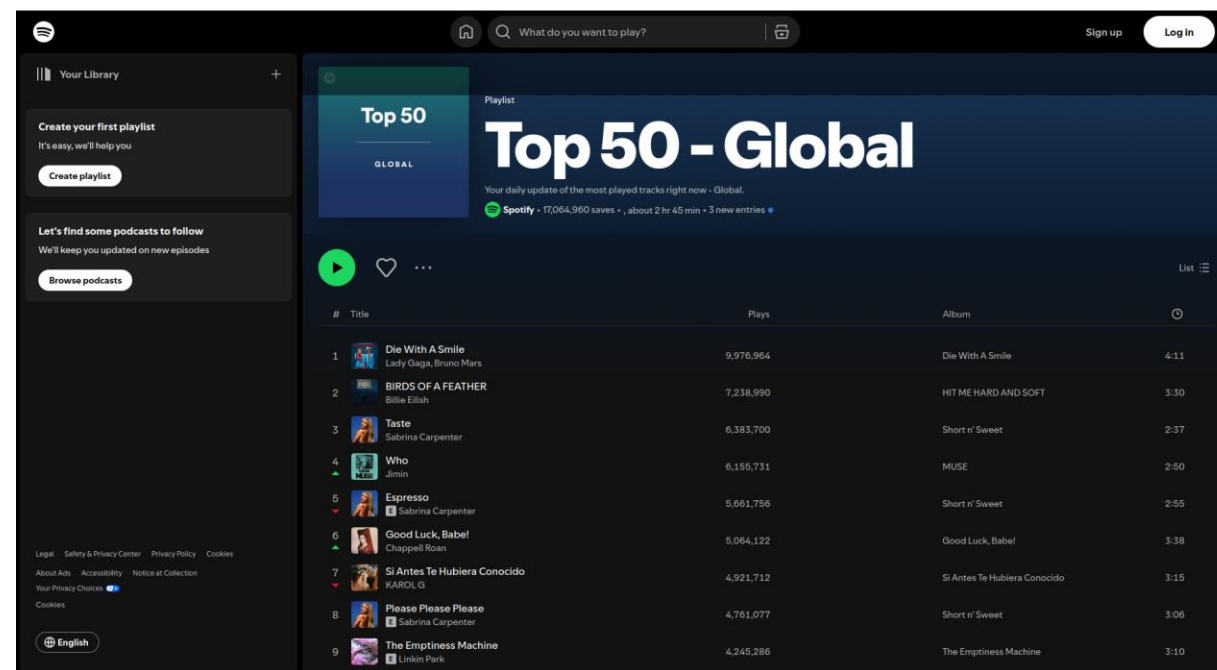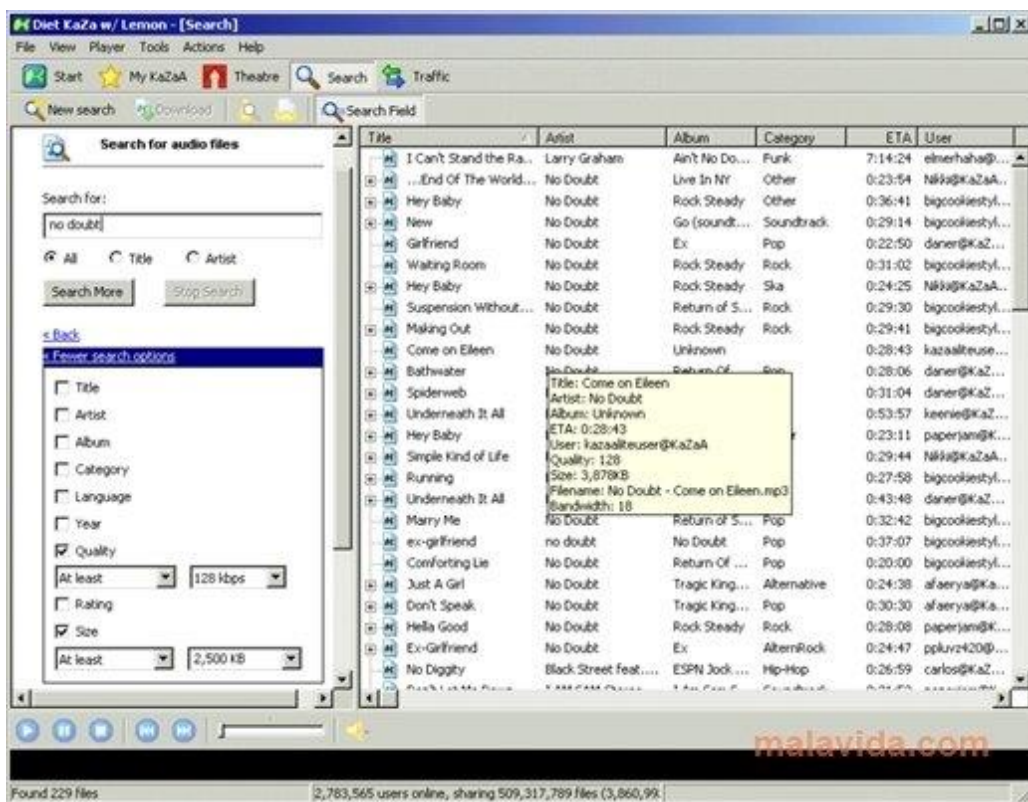
# User expectations are evolving



Kids these days are accustomed to higher-level abstractions over their data.

https://www.theverge.com/22684730/students-file-folder-directory-structure-education-gen-z

# Music Access Has Changed

c. 2000: folders of searchable MP3 files

Now: searchable collections of songs

# Comparison of Data Access Models

| Serves | Scientific Data | Music |
|---|---|---|
| Files, metadata in filenames | SFTP | Kazaa, Napster |
| Files with rich metadata | DataFed, SciCat, Rucio ... | iTunes |
| *Datasets* with rich metadata | ArrayLake, DVID, h5grove, **Tiled** ... | Apple Music, Spotify |

Critical Distinction:

**Can the service open and understand the data it is serving?**

National Synchrotron Light Source II

# File Catalogs Enable...

- User interfaces leveraging rich metadata

- Search

- File-based download

- Cross-site replication

```
https://...?filter=...&select=...
```

# Structured Data Services Additionally Enable...

- Heterogenous data storage, including blob stores key-value stores, databases, detector memory...

- Partial and "chunk-aware" access: Download a region of interest first, or download chunks in parallel

- Transcoding: HDF5 or TIFF -> PNG

```
https://...?slice=50,100:200&format=image/png

https://...?column=temperature&column=intensity&format=text/csv
```

# Important!
# You can still get your data out

- *Apple Music* and *Spotify* don't let you access the underlying storage medium.

- Sometimes scientific applications have good reasons to do this.
    - Performance: Go around the service and straight to GPFS.
    - Backward-compatibility with existing workflows

- Structural data services can still simply provide the filepaths (or database/blob URIs) or serve raw assets to clients

National Synchrotron Light Source II

# A Structured Data Service Knows More

**File Catalog Database Sketch:**

- JSON Metadata blob

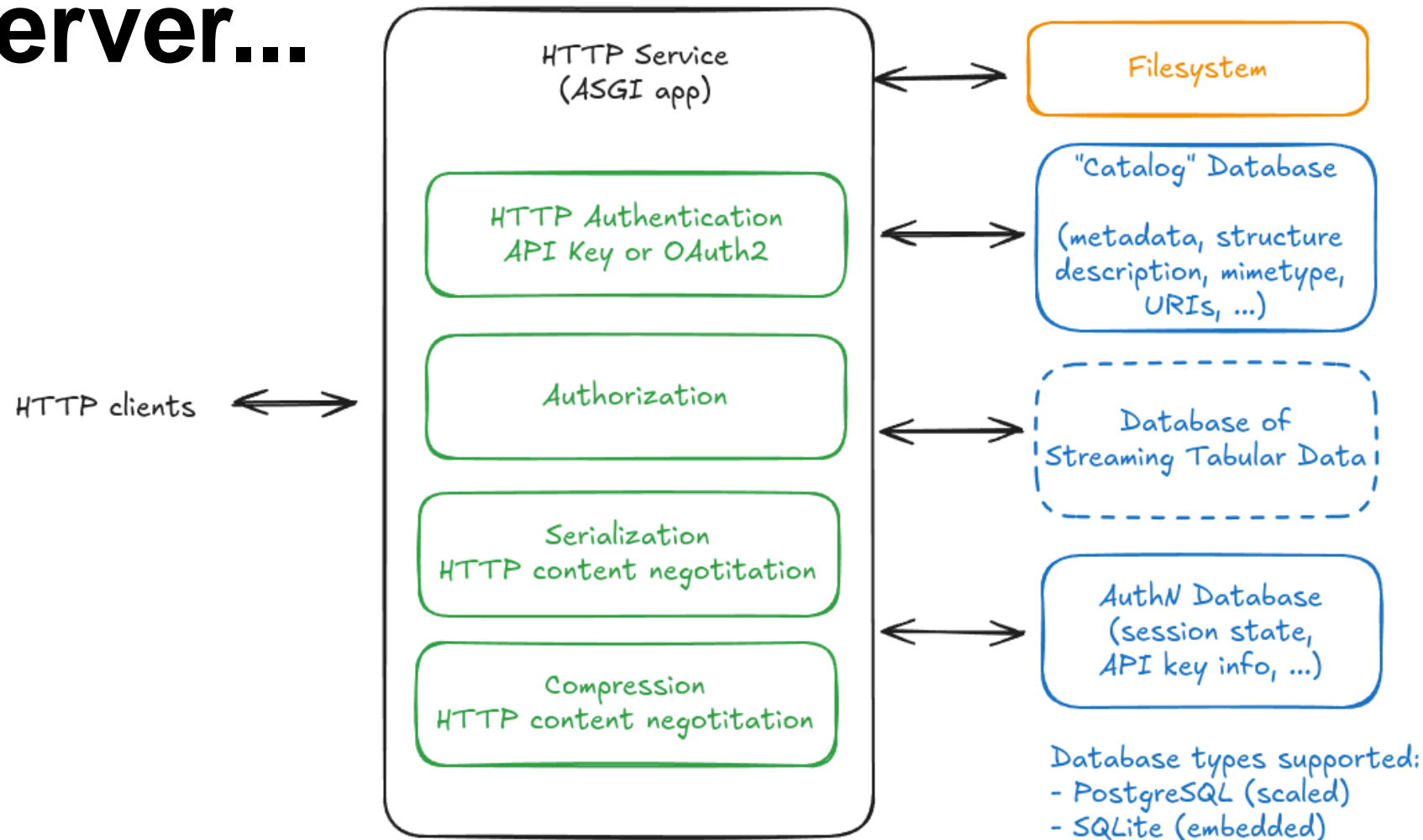- Filepath

**Structured Data Service Database Sketch:**

- JSON Metadata blob

- Data Structure Description
  (array data type and shape,
  table column names and data types,…)

- Mimetype (e.g. 'image/tiff')

- Parameters for accessing data within file
  (e.g. HDF5 dataset path)

- URI(s)

# *Tiled* is...

1. A secure HTTP structural data service, for reading and writing

2. A Python client to that service, like generalized h5py

3. A proof-of-concept React app on that service

… with a budding ecosystem of third-party client applications

# 1. The server...



HTTP Service (ASGI app)

Filesystem

HTTP Authentication
API Key or OAuth2

"Catalog" Database
(metadata, structure description, mimetype, URIs, ...)

Authorization

Database of Streaming Tabular Data

Serialization
HTTP content negotitation

AuthN Database
(session state, API key info, ...)

Compression
HTTP content negotitation

Database types supported:
- PostgreSQL (scaled)
- SQLite (embedded)

HTTP clients

NOBUGS 2024 – Tiled in the Context of Data and Metadata Services – Daniel Allan

# 2. The Python client...

```
>>> from tiled.client import from_uri

>>> client = from_uri("http://localhost:8000")
>>> client
<Container {'some_image', ...} ~500000 entries>

>>> client['some_image'][:]
array([...])

>>> client.write_dataframe({"x": [1, 2, 3], "y": [4, 5, 6]}, key='some_table')
```

# 3. The prototype React app...
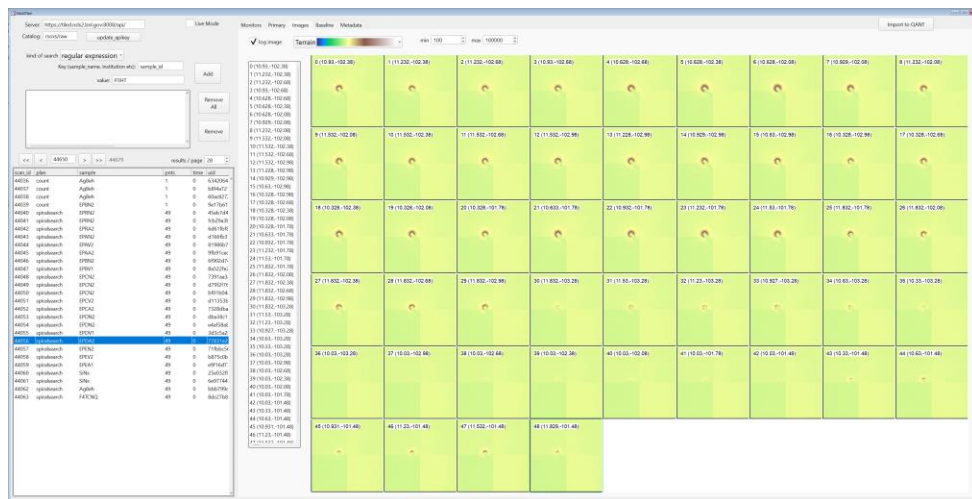
# Third-party client applications...

Eliot Gann's custom (scientist-written!) **Igor** program loads data from Tiled over HTTP, integrating search and viz

**PyMCA** integration is in progress!

And of course it works from `curl`...

# Tiled as an index and an HTTP transport

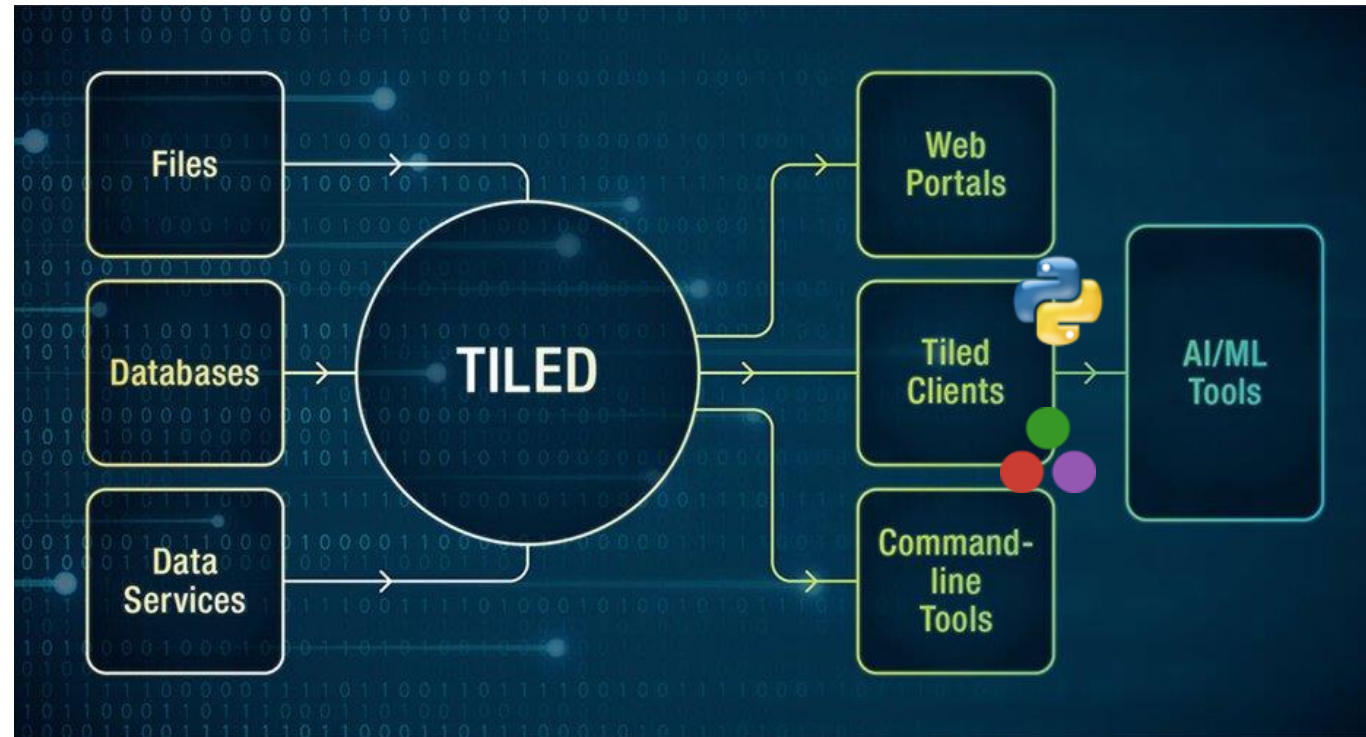| | Index | HTTP transport |
|---|---|---|
| **Reading** | Tiled provides JSON with URI, e.g. `file://.../data.h5` | Download data from Tiled in any supported format |
| **Writing** | Register externally-managed files for Tiled to index and serve | Upload data to Tiled in any supported format for Tiled to store |

# Vital Statistics for *Tiled*

- Developed in the open, **BSD-3** license

- First commit February 2021

- **38** unique code contributors with ~10 unique affiliations

- Under Bluesky **Governance**, may split out on its own someday

- Early usage at: NIST, APS, ALS, BESSY-II, Australian Synchrotron

- **Security** periodically intensively vetted by a third-party cybersecurity penetration testing

# *Tiled* incorporates the input and experience of many people

- **NSLS-II**: Dan Allan, Stuart Campbell, Thomas Caswell, Padraic Shafer, Marcus Hanwell, Juan Marulanda, Kari Barry, Eugene Matviychuk, Seher Karakuzu, Hiran Wijesinghe

- **ALS**: Dylan McReynolds, Joseph Kleinhenz, Wiebke Koepp

- Builds on open-source collaboration with Martin Durant (**Anaconda, Inc.**) with involvement from Garrett Bischof (**NSLS-II**).

- Contributions from the **IRIS-HEP** collaboration to add support for AwkwardArray

# *Tiled:* A Structured API to Data

- **Search** on metadata
- **Slice** into remote datasets
- **Transcode** between formats
- **Download** partial or whole datasets
- Or **find** data storage location(s) for direct access
- Implement web **security** standards and authorization

# Links

Demo: https://tiled-demo.blueskyproject.io/

Documentation: https://blueskyproject.io/tiled

Code: https://github.com/bluesky/tiled

Contact: Daniel Allan <dallan@bnl.gov>