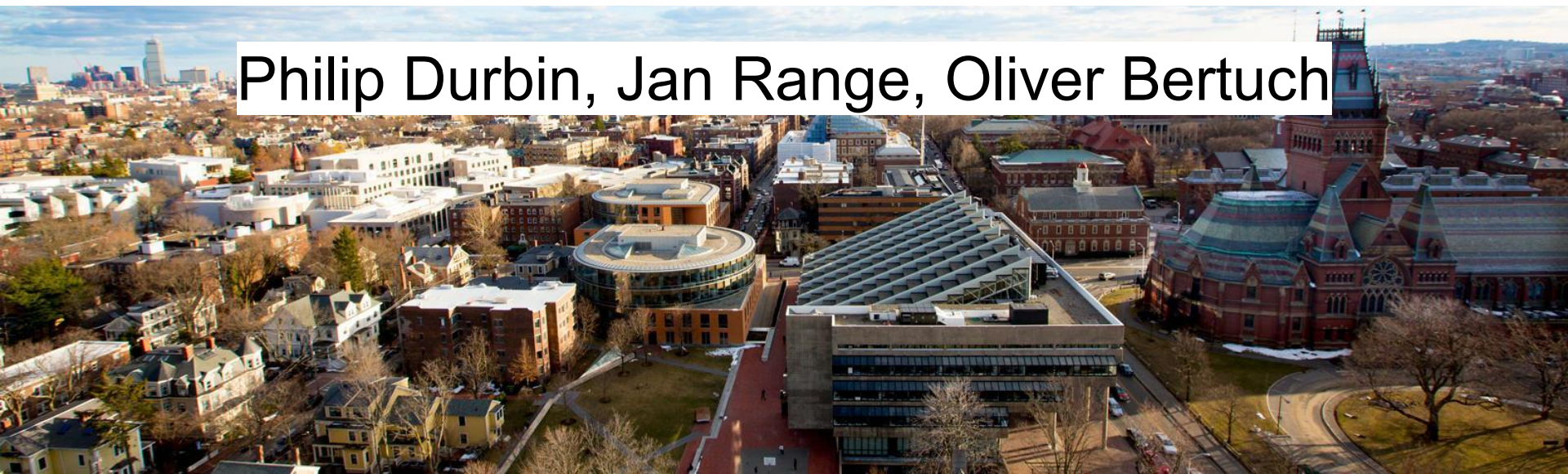# Distributed Metadata and Data with Dataverse

Philip Durbin, Jan Range, Oliver Bertuch

# Agenda

- What is Dataverse?
- Distributed metadata: OAI-PMH (Harvesting)
- Distributed data
- Dataverse and DataLad
- Future architecture for a mixed-style storage configuration in Jülich DATA[1]
- Use of Dataverse APIs to manage data with Python (pyDataverse)

(1): Jülich DATA is the institutional data repository of Forschungszentrum Jülich GmbH, based on Dataverse.

# What is Dataverse?

# What is Dataverse?



116 Installations



The Dataverse® Project

https://dataverse.org



https://doi.org/10.7910/DVN/TJCLKP

# Features

## OAI-PMH (Harvesting)

Gather and expose metadata from and to other systems using standardized metadata formats: Dublin Core, Data Document Initiative (DDI), OpenAIRE, etc. More information.

## Backend storage on S3 or Swift

Choose between filesystem or object storage, configurable per collection and per dataset. More information.

## Direct upload and download for S3

After a permission check, files can pass freely and directly between a client computer and S3. More information.

https://dataverse.org/software-features

# DataverseTV: https://dataverse.org/dataversetv

### Introduction to Data Management and FAIR Data

Katie Mika
Data Services Librarian
Harvard Library & Institute for Quantitative Social Science
katherine_mika@harvard.edu

**Katie Mika**
*2023-09-20*

Introduction to Data Management & FAIR

### Data Handling on The Dataverse Project

Sonia Barbosa, Gustavo Durand, Amber Leahey

**Sonia Barbosa, Gustavo Durand, Amber Leahey**
*2023-10-11*

Introduction to Data Handling on the Dataverse Platform

### Introduction to Metadata in Dataverse Repositories
For Researchers and Support Staff

November 15, 2023
Julian Gautier
IQSS, Harvard University

**Julian Gautier**
*2023-11-15*

Introduction to Metadata in Dataverse Repositories: For Researchers and Support Staff

### Introduction to Dataverse APIs

Gustavo Durand, IQSS, Harvard University
Jim Myers, GDCC

**Gustavo Durand, James Myers**
*2024-01-31*

Introduction to Dataverse APIs

### SESSION 2—MARRIAGE OF DATA AND COMPUTE
### BRINGING THE DATA CLOSE TO THE COMPUTE AT HARVARD DATAVERSE
Afternoon of February 28

**Stefano Iacus**
Director of Data Science and Product Research, IQSS, Harvard

**Stefano Iacus**
*2024-02-28*

Bringing Data Close to Compute at Harvard Dataverse

### Intermediate Curation on the Dataverse Supported Repository

**Sonia Barbosa**
*2023-12-04*

Intermediate Curation on the Dataverse Platform

# Distributed Metadata

# Dataset discoverability

Datasets are made discoverable by a variety of methods.

**Contents:**

- DataCite Integration
- OAI-PMH (Harvesting)
- Machine-Readable Metadata on Dataset Landing Pages
  - Dublin Core HTML Meta Tags
  - Schema.org JSON-LD Metadata
  - Signposting
- Additional Discoverability Through Integrations

https://guides.dataverse.org/en/6.2/admin/discoverability.html

# OAI-PMH (Harvesting)

## The Open Archives Initiative Protocol for Metadata Harvesting

Protocol Version 2.0 of 2002-06-14
Document Version 2015-01-08
http://www.openarchives.org/OAI/2.0/openarchivesprotocol.htm

Previous protocol version: **Protocol Version 1.1 of 2001-07-02**
**Instructions** for migrating from Version 1.1 to 2.0
**Implementation Guidelines**

### Editors

The OAI Executive:
Carl Lagoze <lagoze@cs.cornell.edu > -- Cornell University - Computer Science
Herbert Van de Sompel <herbertv@lanl.gov > -- Los Alamos National Laboratory - Research Library

From the OAI Technical Committee:
Michael Nelson <m.l.nelson@larc.nasa.gov > -- NASA - Langley Research Center
Simeon Warner <simeon@cs.cornell.edu > -- Cornell University - Computer Science

http://www.openarchives.org/OAI/2.0/openarchivesprotocol.htm

# Harvesting clients

**Harvesting Clients**

29
Clients

82186
Datasets

Manage Clients

Harvard Dataverse > Dashboard >

ℹ️ **Manage Harvesting Clients** – Harvesting can be scheduled to run at a specific time or on demand. Harvesting can be initiated here or via the REST API. ✕

+ Add Client

| Nickname/URL | Dataverse | Last Run | Last Result | Actions |
|---|---|---|---|---|
| acss_historic_shapefiles<br>https://dataverse.theacss.org/oai | Historic shapefiles | Fri Mar 22 04:00:00 EDT 2024 | SUCCESS; 0 harvested, 0 deleted, 0 failed. | ▶ ✏ 🗑 |
| borealis<br>https://borealisdata.ca/oai | Borealis Harvested Dataverse | Thu Feb 29 16:13:26 EST 2024 | SUCCESS; 19207 harvested, 0 deleted, 453 failed. | ▶ ✏ 🗑 |
| cifor<br>https://data.cifor.org/oai | CIFOR Harvested Dataverse | Sun Mar 17 05:00:00 EDT 2024 | SUCCESS; 0 harvested, 0 deleted, 15 failed. | ▶ ✏ 🗑 |
| cimmyt<br>https://data.cimmyt.org/oai | CIMMYT Harvested Dataverse | Sun Mar 17 23:00:00 EDT 2024 | SUCCESS; 0 harvested, 0 deleted, 405 failed. | ▶ ✏ 🗑 |
| cimmyt-iwyp<br>https://data.cimmyt.org/oai | CIMMYT Harvested Dataverse | Sat Mar 16 23:00:00 EDT 2024 | SUCCESS; 0 harvested, 0 deleted, 0 failed. | ▶ ✏ 🗑 |
| cimmyt-software<br>https://data.cimmyt.org/oai | CIMMYT Harvested Dataverse | Sun Mar 17 07:00:00 EDT 2024 | SUCCESS; 0 harvested, 0 deleted, 17 failed. | ▶ ✏ 🗑 |
| cip<br>https://data.cipotato.org/oai | International Potato Center Harvested Dataverse | Mon Mar 04 04:00:01 EST 2024 | FAILED | ▶ ✏ 🗑 |
| dataverse-nl<br>https://dataverse.nl/oai | DataverseNL Harvested Dataverse | Mon Mar 04 03:44:20 EST 2024 | SUCCESS; 18 harvested, 0 deleted, 0 failed. | ▶ ✏ 🗑 |

# Harvesting formats (Dublin Core, DDI, DataCite, etc.)

# Harvesting server (sets)

**Harvesting Server**

| | |
|---|---|
| Enabled | 33 |
| Status | Sets |

Manage Server

Harvard Dataverse > Dashboard >

ℹ **Manage Harvesting Server** – Define sets of local datasets that will be available for harvesting by remote clients.          ✕

**OAI Server**   [ Enabled ▼ ]

➕ Add Set

| OAI setSpec/Description | Definition Query | Datasets | Actions |
|---|---|---|---|
| **DEFAULT**<br>The default, "no name" set. The OAI server will serve the records from this set when no "setspec" argument is specified by the client. | All Published Local Datasets | 91,283 datasets (91,283 records exported, 12,072 marked as deleted) | ▶ ✏ 🗑 |
| **AfricaRice**<br>For Thompson/Reuters | subtreePaths:"/1531" | 124 datasets (124 records exported, 84 marked as deleted) | ▶ ✏ 🗑 |
| **AllianceBioversityCIAT**<br>Datasets published in the The Alliance of Bioversity International and CIAT Dataverse at https://dataverse.harvard.edu/dataverse/AllianceBioversityCIAT | subtreePaths:"/3538956" | 188 datasets (188 records exported, 1 marked as deleted) | ▶ ✏ 🗑 |
| **Bioversity**<br>Datasets published in the Bioversity International Dataverse (https://dataverse.harvard.edu/dataverse/Bioversity) | subtreePaths:"/1823" | 108 datasets (108 records exported, 0 marked as deleted) | ▶ ✏ 🗑 |

# More results through harvesting



distributed | Advanced Search | + Add Data

- ✔ 👤 **Dataverses (60)**
- ✔ 📄 **Datasets (4,455)**
- ✔ 📊 **Files (1,975)**

**1 to 10 of 6,490 Results** | ⇅ Sort ▾

**Dataverse Category**
Research Project (31)
Organization or Institution (10)
Researcher (10)
Research Group (3)
Department (1)

The **Distributed** AI Research Institute  (DAIR Institute)
Feb 14, 2024

The Legacy of Spatial Apartheid Dataset

**Metadata Source**
Harvard Dataverse (3,419)
Harvested (3,071)

---

distributed | Advanced Search | + Add Data

- ✔ 👤 **Dataverses (0)**
- ✔ 📄 **Datasets (2,730)**
- ✔ 📊 **Files (341)**

Metadata Source: Harvested ✖

**1 to 10 of 3,071 Results** | ⇅ Sort ▾

**Publication Year**
2015 (1,117)
2011 (350)
2007 (37)

**Subject**
Agricultural Sciences (331)
Social Sciences (150)
Earth and Environmental Sciences (145)
Arts and Humanities (54)
Medicine, Health and Life Sciences (38)

20210126 - Ocean Biology **Distributed** Active Archive Center - CoreTrustSeal Requirements 2020-2022
Jan 23, 2024 - DataverseNL Harvested Dataverse

Ocean Biology Distributed Active Archive Center, 2024, "20210126 - Ocean Biology Distributed Active Archive Center - CoreTrustSeal Requirements 2020-2022", https://doi.org/10.34894/7N4H2H

CoreTrustSeal certification
Author Name: Ocean Biology Distributed Active Archive Center
This Dataset is harvested from our partners. Clicking the link will take you directly to the archival source of the data.

20210126 - Physical Oceanography **Distributed** Active Archive Center - CoreTrustSeal Requirements 2017-2019
Jan 23, 2024 - DataverseNL Harvested Dataverse

Physical Oceanography Distributed Active Archive Center, 2024, "20210126 - Physical Oceanography Distributed Active Archive Center - CoreTrustSeal Requirements 2017-2019", https://doi.org/10.34894/PYZQCJ

# 2002? Anything newer?

Defining the next generation repository

**NEXT GENERATION REPOSITORIES**

**COAR**
Confederation of Open Access Repositories

Vision & Objectives    Principles    Design Assumptions    User Stories    Behaviours

Technologies    About

## ResourceSync

ResourceSync is a specification based on Sitemaps that can be used by repository managers to provide information that allows third-party systems to remain in sync with the resources in their repository as they evolve, i.e. are created, updated, deleted. Whereas basic Sitemaps allow exposing a repository inventory and crawl-related metadata, ResourceSync adds ways to expose changes only, and to provide expressive synchronization-related metadata as well as typed links for further discovery. ResourceSync can be used for discovery and synchronization of both content and metadata and uses the Sitemaps XML format. http://www.openarchives.org/rs/toc

### Related behaviours

- Resource Transfer
- Batch Discovery
- Collecting and Exposing Activities

https://ngr.coar-repositories.org/technology/resourcesync/

---

**You**
Are there any newer alternatives to OAI-PMH?

**ChatGPT**
As of my last update in January 2022, the Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH) has been a widely adopted standard for harvesting metadata from repositories. However, the landscape of digital repositories and metadata harvesting is continually evolving, and there may be newer alternatives or complementary protocols that have emerged since then.

Some potential alternatives or related protocols that have gained attention in the scholarly communication and digital repository communities include:

1. **ResourceSync**: This is a synchronization framework for the web that is more flexible than OAI-PMH. It allows for more granular synchronization of resources and can be used for a variety of purposes beyond metadata harvesting.
2. **IIIF (International Image Interoperability Framework)**: While primarily focused on image interoperability, IIIF also supports metadata and could potentially be used in conjunction with other protocols for harvesting and sharing metadata.
3. **JSON-based APIs**: Many modern repositories and platforms are moving towards RESTful APIs that return data in JSON format. While not a replacement for OAI-PMH per se, these APIs can be used for harvesting metadata and accessing repository content.
4. **GraphQL**: This query language for APIs offers a more flexible and efficient way to retrieve data compared to traditional RESTful APIs. Some repositories and platforms are starting to offer GraphQL APIs for accessing and harvesting metadata.
5. **Linked Data and RDF**: The Semantic Web technologies, including Linked Data and RDF (Resource Description Framework), offer ways to publish and link data in a standardized and interoperable manner. Some repositories are exploring these technologies for exposing metadata and content.
6. **Webmention and WebSub**: These are protocols for web-based notifications and updates, which can be used to notify harvesting services of changes to content or metadata.

It's important to note that the choice of protocol or technology will depend on the specific requirements and context of the digital repository or system in question. Some repositories may continue to use OAI-PMH due to its widespread adoption and established ecosystem, while others may opt for newer alternatives or a combination of different protocols and technologies.

If you're interested in exploring alternatives to OAI-PMH, I would recommend researching the latest developments in digital repository technologies, attending relevant conferences and workshops, and consulting with experts in the field to determine the best approach for your specific needs.

# Distributed Data

- File Storage
  - Multi-store Basics
  - Labels for File Stores
  - File Storage
  - Swift Storage
    - Setting up Compute with Swift
  - Amazon S3 Storage (or Compatible)
    - First: Set Up Accounts and Access Credentials
      - Preparation When Using Amazon's S3 Service
      - Preparation When Using Custom S3-Compatible Service
      - Manually Set Up Credentials File
      - Console Commands to Set Up Access Configuration
    - Second: Configure Your Dataverse Installation to use S3 Storage
      - List of S3 Storage Options
      - Credentials via MicroProfile Config
      - Reported Working S3-Compatible Storage
      - Migrating from Local Storage to S3
  - Trusted Remote Storage
  - Globus Storage
  - Temporary Upload File Storage

https://guides.dataverse.org/en/6.2/installation/config.html#file-storage

⬇ Publication Guidelines

Data ▾     About ▾     Get Help ▾     For Depositors ▾     🔍

## X-Ray Diffraction data from Lin28A/let-7g microRNA complex, source of 3TS2 structure

VISUALIZE IN 3D >>

Data DOI: 10.15785/SBGRID/1 | ID: 1
Publication DOI: 10.1016/j.cell.2011.10.020
3TS2 Coordinates: Viewer, PDB (RCSB) (PDBe), MMDB
Sliz Laboratory, Harvard Medical School
Release Date: April 10, 2015

### Data Access Instructions

1. If this dataset is locally available, it should be accessible at `/programs/datagrid/1`

2. To download this dataset, please run the following command from your Terminal on a Linux or OS X workstation:

`'rsync –av rsync://data.sbgrid.org/10.15785/SBGRID/1 .'` (Harvard Medical School, USA)

Depending on your location, faster access may be available from a Tier 1 site closer to your location

`'rsync –av rsync://sbgrid.icm.uu.se/10.15785/SBGRID/1 .'` (Uppsala University, Sweden)

`'rsync –av rsync://sbgrid.pasteur.edu.uy/10.15785/SBGRID/1 .'` (Institut Pasteur de Montevideo, Uruguay)

`'rsync –av rsync://sbgrid.ncpss.org/10.15785/SBGRID/1 .'` (Shanghai Institutes for Biological Sciences, China)

3. After the transfer is completed, please issue the following command to verify data integrity:

`'cd 1 ; shasum –c files.sha'`

Storage requirements: 1.6G

https://data.sbgrid.org/dataset/1

# Dreams of distributed data

"An idea of: synchronizing the data between the Harvard Dataverse and the soon-to-be-up [Foobar Dataverse] (and possibly other major locations) to make a connected dataverse. I believe that it would enhance the usability of the system for the global research community and possibly distribute the traffic better."

# Amazon CloudFront



before

after

https://aws.amazon.com/blogs/networking-and-content-delivery/amazon-s3-amazon-cloudfront-a-match-made-in-the-cloud/

# How It Works

**Storj Nodes** — Supply

Thousands of shared hard drives store pieces of data on the network, without access to any complete file or usable data. Node operators fairly (and profitably) compensated.

**Applications** — Demand

Client applications store encrypted and encoded files split into fragments and stored across the distributed storage network.

**Satellites** — Network

The Storj network enables applications to store data, ensures data reliability, manages access controls, and pays storage nodes.

Satellites

Uplink

S3

Customer Application

Storage Node Network

https://osf.io/3txy4
https://www.youtube.com/watch?v=3ek7F_Dxcjk&t=2102s

# Dataverse and DataLad



https://docs.datalad.org/projects/dataverse



After my @dataverseorg talk at #FOSDEM @eknahm and @yarikoptic from @datalad offered me and @Waaaloo beers and great conversation. Let's integrate!

10:13 AM · Feb 2, 2020 from Brussels, Belgium

# Future architecture for a mixed-style storage configuration in Jülich DATA

# What is Jülich DATA (now)?



- Status quo: institutional data registry since 2020
- Hosting Central Library, Open Science Team for FZJ (~2.5k researchers)
- Zero storage policy / ideal

# Publication personas: current use cases

## Alice

Uploads to and publishes in external repository

Registers in Jülich DATA with added institutional metadata

Main reason: reporting for €€€ from Helmholtz

## Bobby

Metadata-only publi-cation in Jülich DATA

Data is referenced only for valid reasons (privacy, size, policy) using a URI to storage location

## Devin

In denial, does not publish anything

¯\_(ツ)_/¯

# Lessons learned

- Only 150 (published) datasets?!
- Metadata input twice (publication and registration) makes users unhappy

- Open Science Team is well established addressee for RDM questions
- We need to turn people away looking for publication options, especially >300GB and no/small community repos
- Chicken & Egg: "No easy solution available? Build our own or don't publish."

- Discoverability of files for metadata-only publications is zero
- FZJ joined SF-DORA in 2023, we want software publications, too

# Infrastructure Updates

# Infrastructure Update Challenges

- Kubernetes platform
- Built in 2020
- Optimized for registry

- How much storage?
- Structured interviews w/ institute data stewards

Status Quo → New platform → New storage → Next up: Features

- Terraform new K8s
- Migrate to contributed upstream containers
- Contribute K8s operator for management

- T.B.D. by autumn 2024
- Onto more adventures!

# Feature Updates

# Challenges to support Alice (Registration)

# Challenges to support Bobby (*mostly* Metadata Publication)



Dataverse support for more file level metadata?

Dataverse integration with Onedata to support data spaces?

Dataverse support for file and dataset relations?

Sounds like 2025…

Development +Upstreaming necessary for better *file discoverability*

Dataverse support for auxiliary file quotas

Dataverse special storage driver for DataLad and git-annex?

Dataverse support for RO-Crate like dataset nesting?

# Publication personas: adding a third use case

## Alice

Uploads to and publishes in community repository

Registers in Jülich DATA

## Bobby

Metadata publication in Jülich DATA

(Most) data is referenced only for valid reasons (privacy, size, policy)

## Charlie

Uploads to and publishes using Jülich DATA

## Devin

In denial, does not publish anything

¯\_(ツ)_/¯

# Challenges for Charlie (Publication)

Sounds like 2025+

**Social aspects:**
1. Free space is always used
2. Good metadata is hard work

More embed. Data Stewards

Apply Quotas

Training Training Training

Metrics?

(Usage, Quality, Amount of work, …)

Computation next to data?

Control institute archive storage?

Storage tiering with tape libs?

Globus integration? (€€€€€€)

University of Stuttgart
Cluster of Excellence „Data-integrated Simulation Sciences"

SimTech

**Use of Dataverse APIs
to manage data with Python**

Jan Range, Philip Durbin, Oliver Bertuch | Distribits| 05-04-2024

# About PyDataverse
## Overview

## Python library to interface Dataverse

- Developed by Stefan Kasberger (until ~2021)

- Implements a subset of essential Dataverse endpoints

## Popular among Datanauts

- 59 stars on GitHub (best on GDCC)

- 72 repositories utilize pyDataverse

# About PyDataverse
## Overview

## Low-level API for Dataverse

- Direct implementation of endpoints

- Maximum control and flexibility

- Supports async requests

## Essentials for Datanauts

- Global management and infos

- CRUD operations for any element

**Collections**

**Datasets**

**Files**

**Info**

**Schemes**

**Roles**

# Caveat of PyDataverse
## Low-level to the core

## Requires deeper knowledge

- Understanding of the native API

- Combination of multiple endpoints for high-level concepts (e.g. Dataset handling)

## Example: Dataset Creation

- Requires two endpoints to be called

- Complex Dataverse JSON to add metadata

```json
{
  "metadataBlocks": {
    "citation": {
      "displayName": "Citation Metadata",
      "name": "citation",
      "fields": [
        {
          "typeName": "title",
          "multiple": false,
          "typeClass": "primitive",
          "value": "My dataset"
        },
        {
          "typeName": "author",
          "multiple": true,
          "typeClass": "compound",
          "value": [
            {
              "authorName": {
                "typeName": "authorName",
                "multiple": false,
                "typeClass": "primitive",
                "value": "John Doe"
              }
            }
          ]
        }
      ]}]}}}
```

# Caveat of PyDataverse
Low-level to the core

## Requires deeper knowledge

- Understanding of the native API

- Combination of multiple endpoints for high-level concepts (e.g. Dataset handling)

## Example: Dataset Creation

- Requires two endpoints to be called

- Complex Dataverse JSON to add metadata



```
"name": "citation",
"fields": [
  {
    "typeName": "title",            Key
    "multiple": false,
    "typeClass": "primitive",
    "value": "My dataset"           Value
  },
  {
    "typeName": "author",
    "mul
    "ty        The User's Interest
    "value": [
      {
        "authorName": {
                        Key    "typeName": "authorName",
                               "multiple": false,
                               "typeClass": "primitive",
                        Value  "value": "John Doe"
        }
      }
    ]
  }
}
```
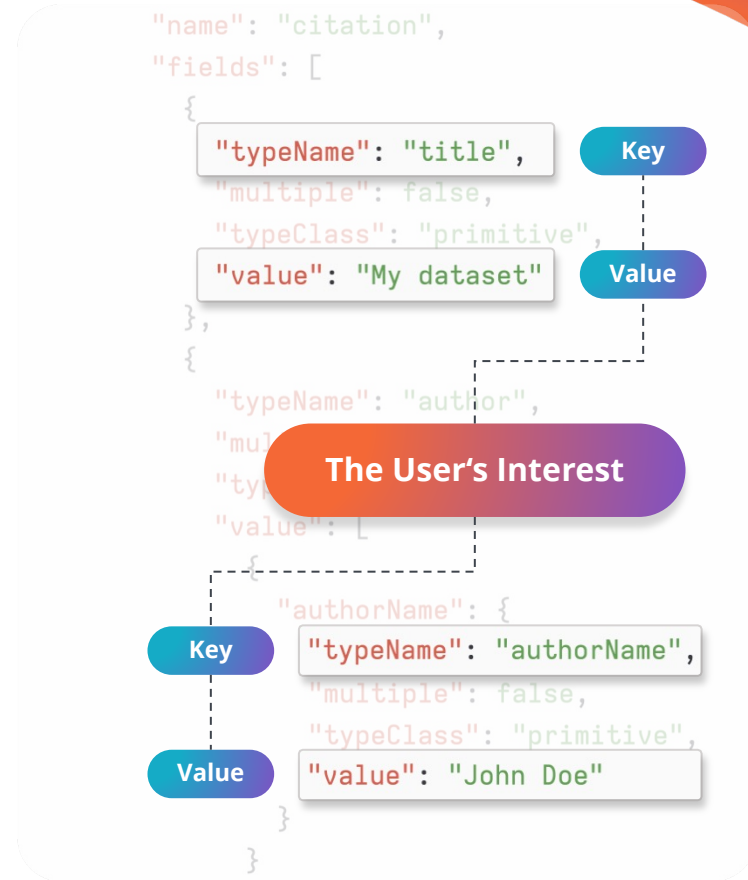
# High-level concepts
## Dynamic metadata objects

- Fetches metadata schemes on connection

- Transforms schemes into Python classes

- Stores technical details in attribute metadata

## High-level recipe

- Combines multiple endpoints into a simple interface

- Improves user-friendliness and convenience

- Built-in type and schema validation

**Metadata schemes**

```python
# Set Citation metadata
dataset.citation.title = 'My dataset'
dataset.subject = ['Other']
dataset.citation.add_author(
    name='John Doe',
    affiliation='My University'
)
```

**Dataverse JSON**

## High-level concepts
## Parallelized upload to S3

- Facilitates direct upload to an S3 bucket

- Complex procedure of multiple endpoints

- Avoids duplicate files by checksum lookup

### High-level recipe

- Abstracts procedure into a simple interface

- Improves user-friendliness and usability

- Asynchronous process to maximize performance



```
(dvuploader) → examples git:(async-requests) ✗ python run.py
```

**Practical example**

Dataset handling with S3

**1**  **Create dataset at demo.dataverse.org**

- Using metadata block objects

- Uploading multiple files to S3

**2**  **Update the dataset and add more files**

→ Google Colab Notebook

# Join us!
## PyDataverse Working Group

**We are open to anyone who wants to participate and contribute!**

→ Bi-weekly meetings

→ Recordings, notes and upcoming dates available at https://py.gdcc.io

💬 *#python* on Zulip → Here