# Clowder

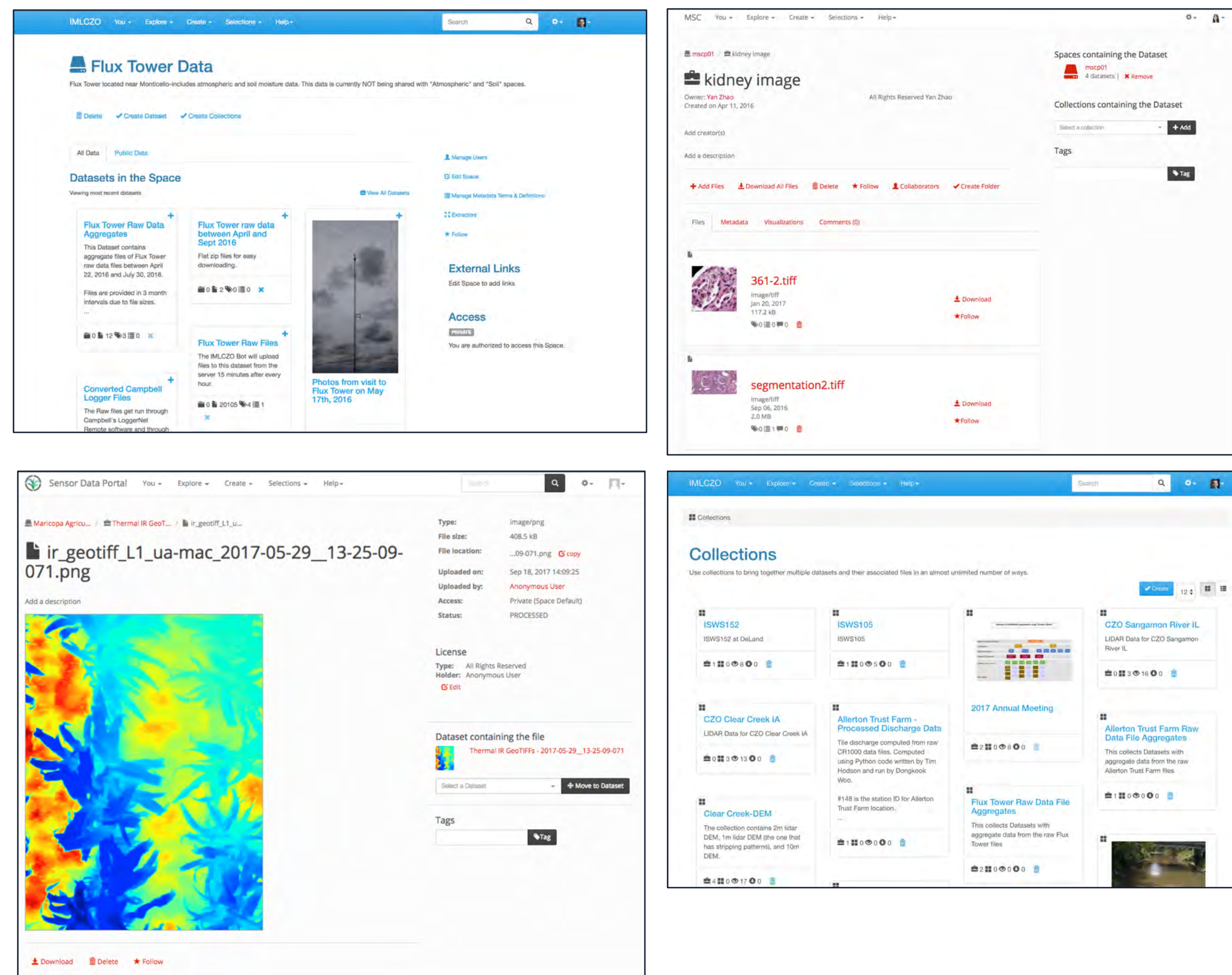# Open Source Data Management for Long Tail Data

Luigi Marini, Rob Kooper, Indira Gutierrez, Max Burnette, Sandeep Puthanveetil Satheesan, Bing Zhang, Mike Lambert, Todd Nicholson, Yan Zhao, Jong Lee, Kenton McHenry
National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign

https://clowder.ncsa.illinois.edu

## Introduction

- Clowder is a customizable and scalable data management system you can install anywhere.
- You can install Clowder in the cloud, on your hardware, or you can partner with NCSA for a custom instance.
- You can contribute to the core software or by creating new metadata extractors and data visualizations.
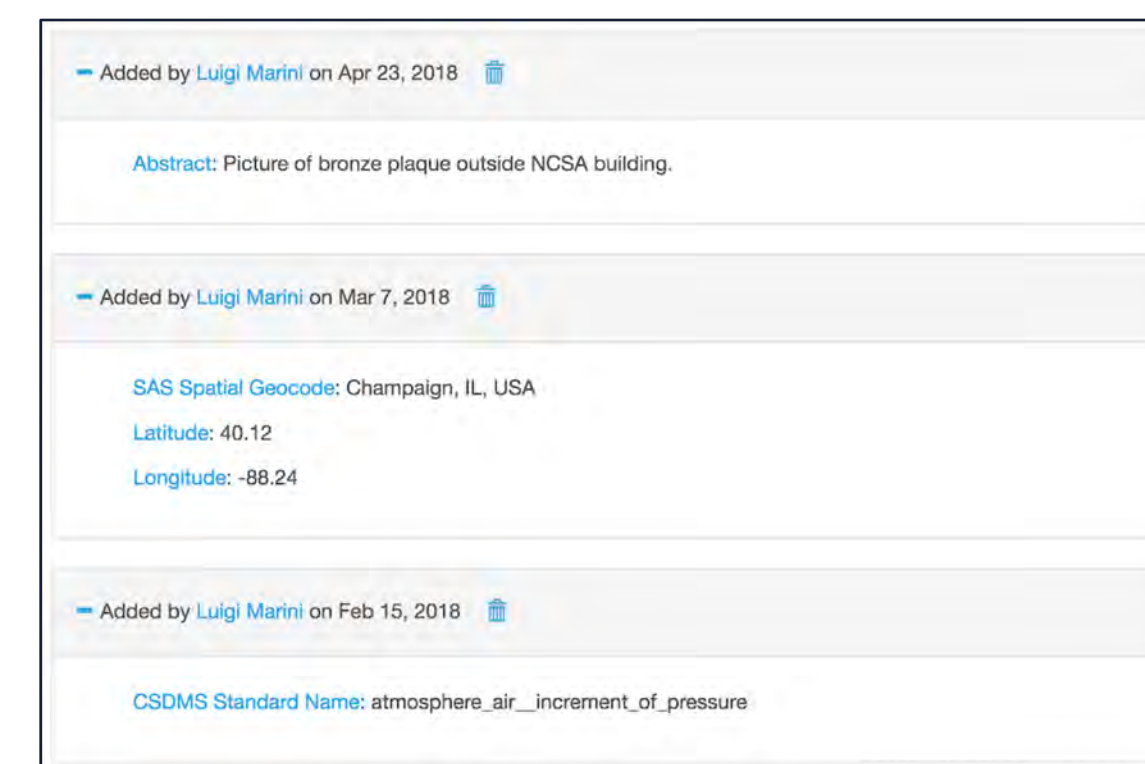
## Organize Data



*Spaces (top-left), Datasets (top-right), Files (bottom-left), Collections (bottom-right).*

## Open Source Community



Used across many different research areas.

## Flexible Metadata Representation

Support for both user-defined and machine-defined metadata. System accepts metadata in a flexible representation based on JSON Linked Data.



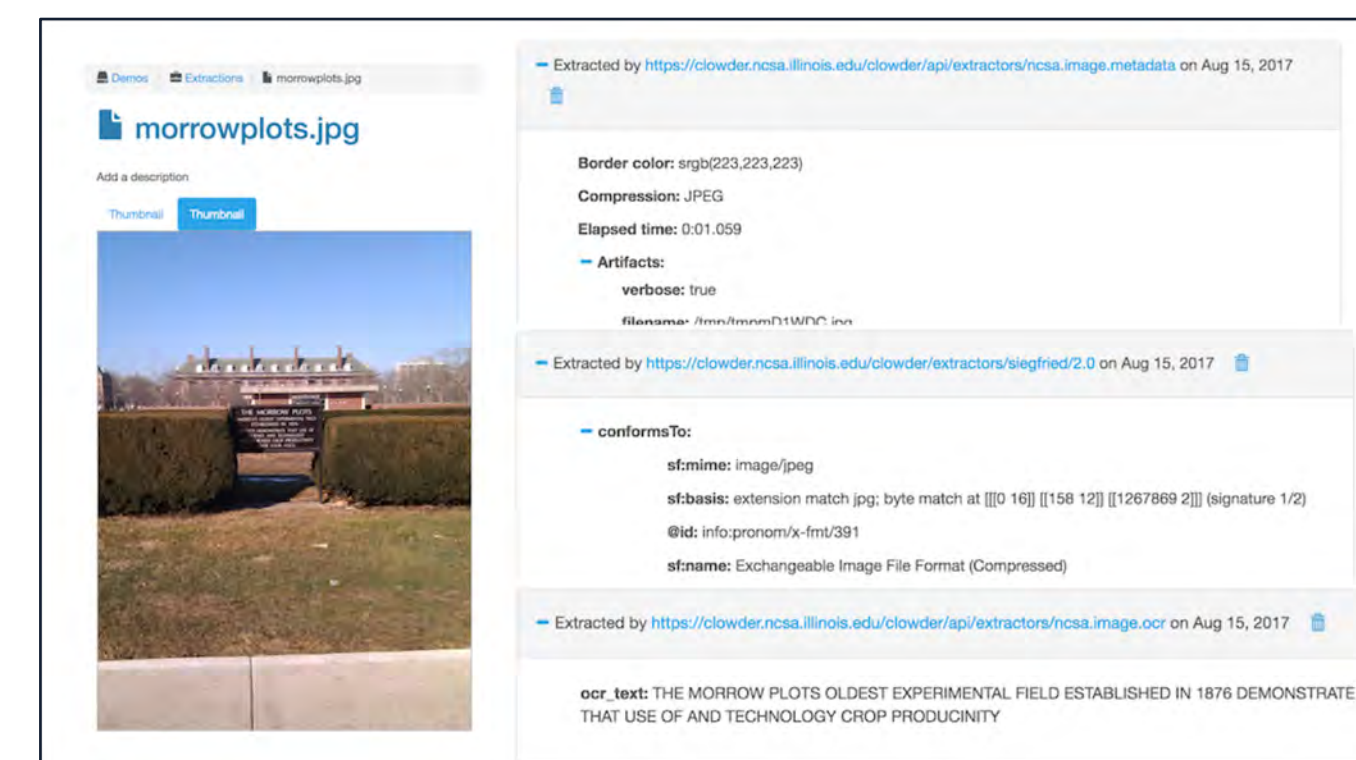*Users can add metadata entries directly from the user interface.*

*Extractors and external clients can attach metadata to files and datasets using the Web service API.*

## Automatic Metadata Extraction

Extend the system by creating new extractors to analyze data. For example, given the image on the right 3 separate extractors computed:
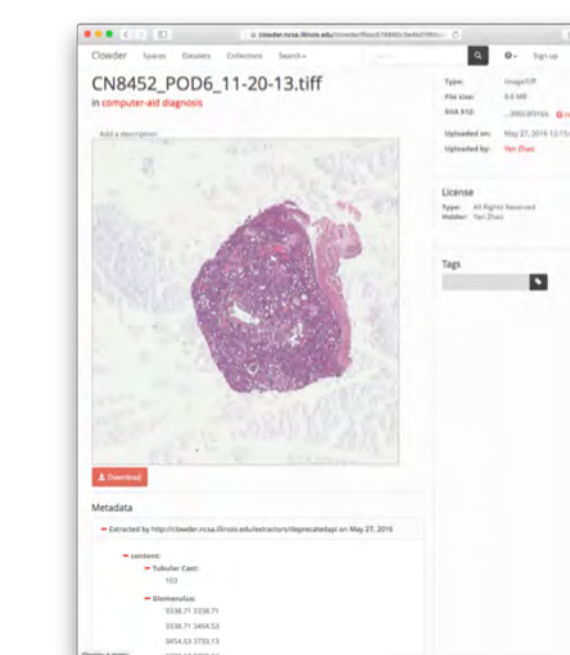- EXIF from image bytes
- Siegfried file format identification
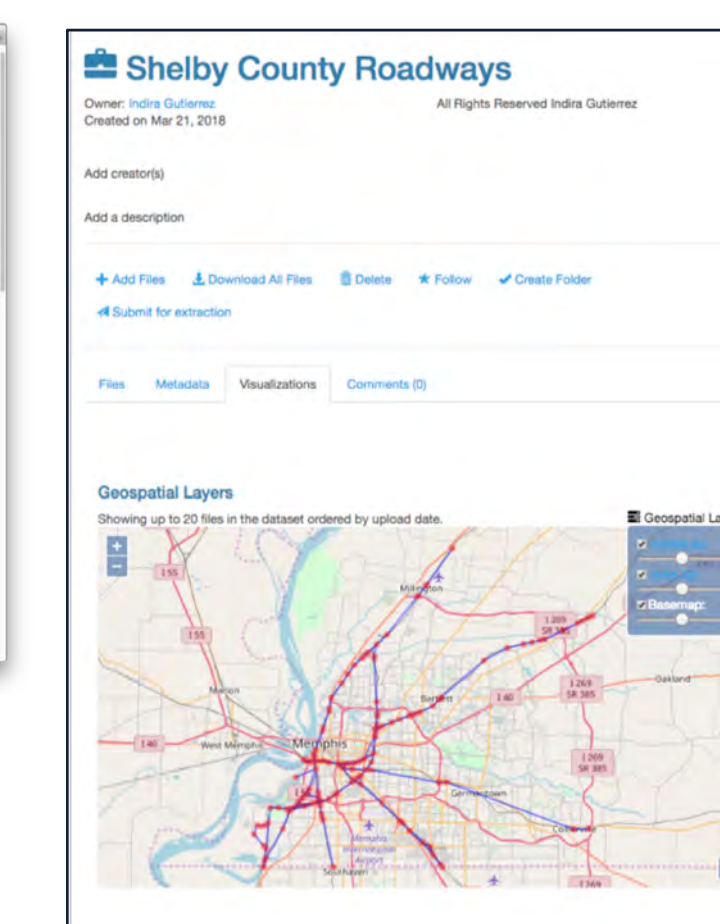- OCR of text in image



Write your extractors using the Clowder Python SDK or in most other programming languages.

## Search

Flexible text search over basic attributes of files, datasets and collections. Advanced search over metadata fields.
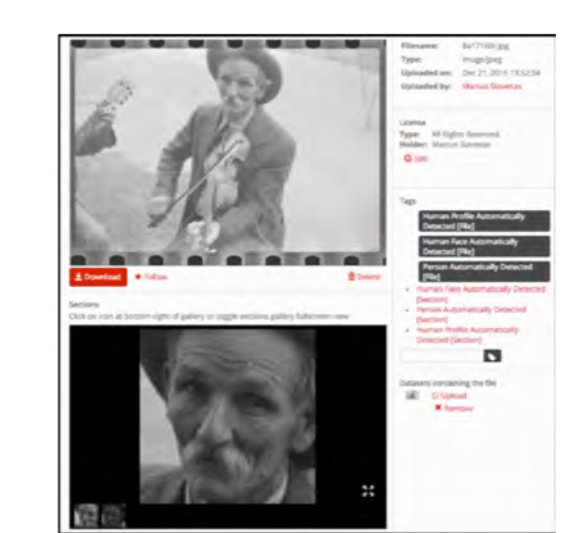


## Data Visualizations

- To preview the content of large files and visualize the information contained in files and datasets in a meaningful way, Clowder provides ways to write Javascript-based visualizations.
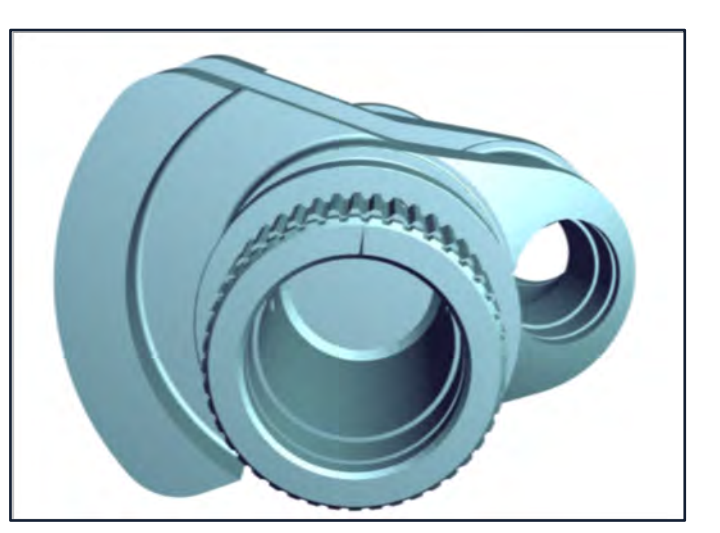- Often these data previews are added by automatic extractions.
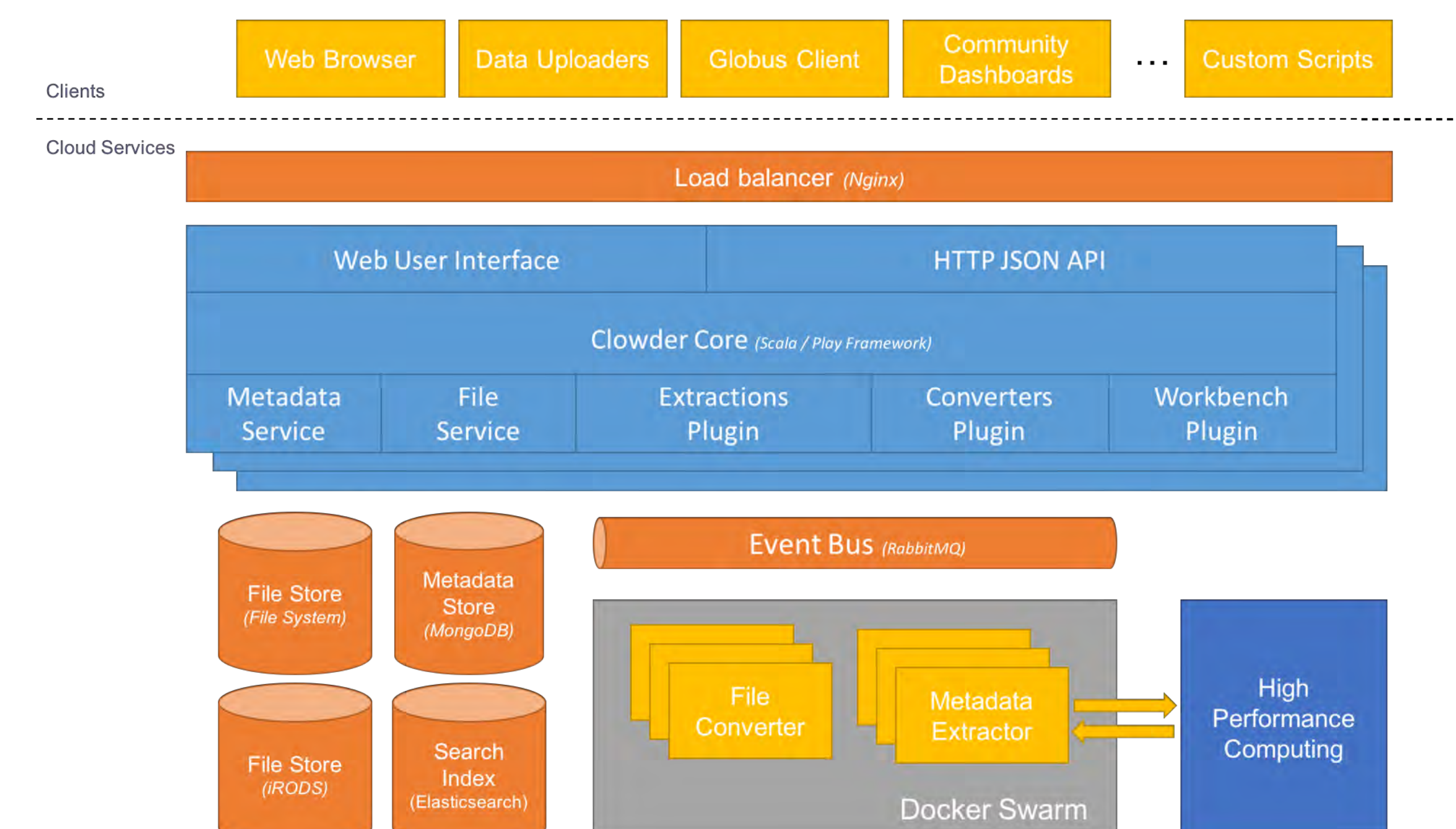


*Cell identification*

*Face identification*

*Interactive 3D object viewer*
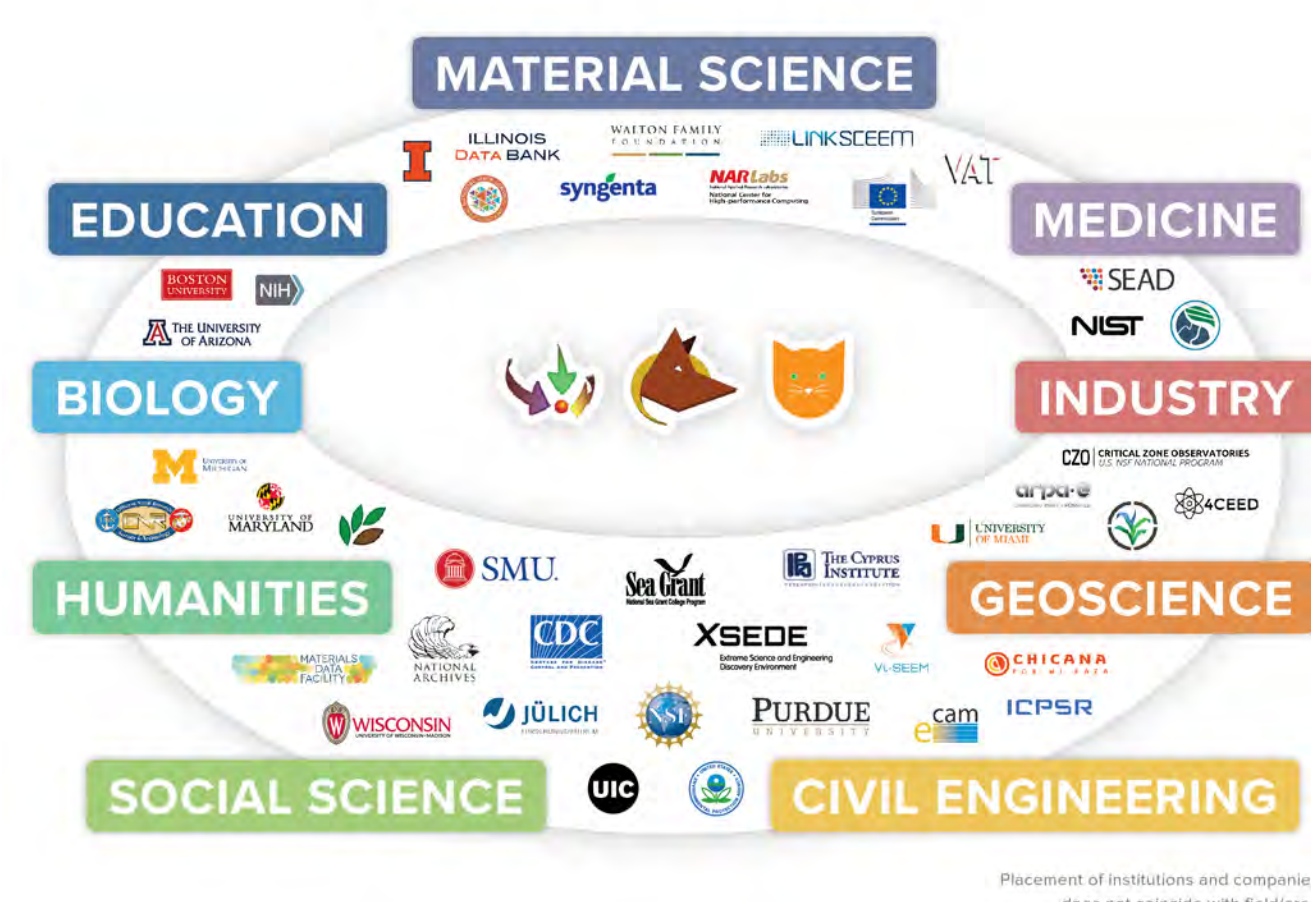
*GIS Shapefile previewer*

## Architecture

Clowder is built from the ground up to be extensible and support new use cases and domains.



- The scalability of the Clowder system has been proven with one instance (TERRA-REF) having 1 PB of data and close to 40 million files.
- The Brown Dog project runs 60 different extractors as services in a docker swarm, which consists of 30 machines. Extractors are scaled elastically based on the number of requests and the swarm runs around 150 instances of these extractors.
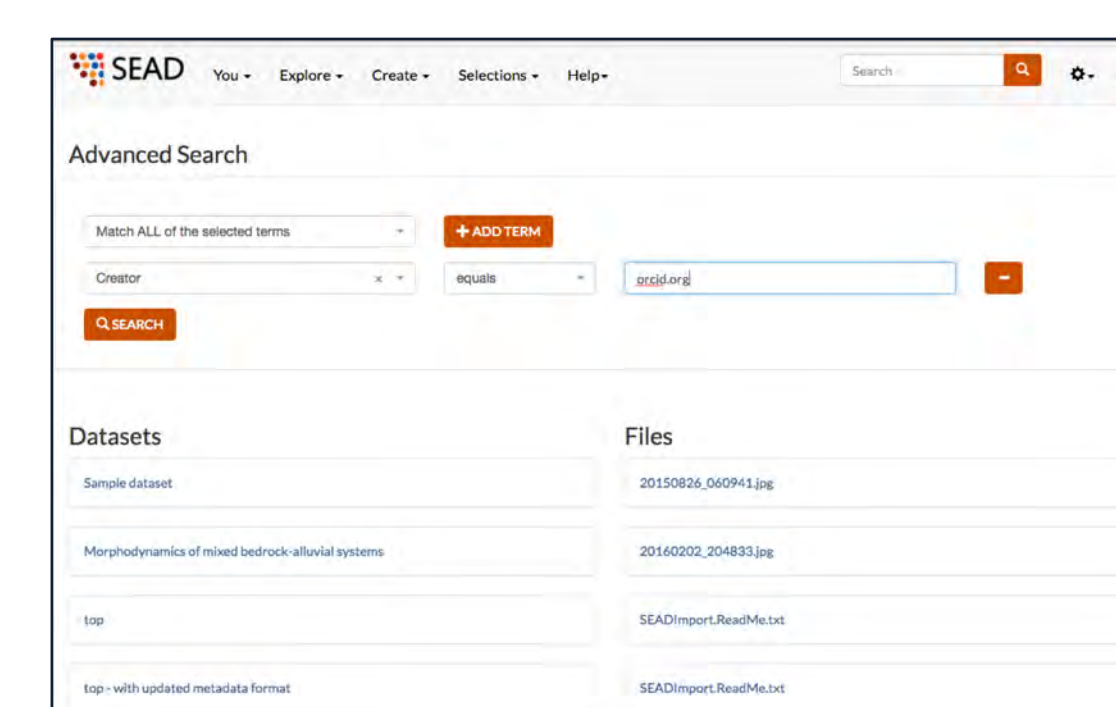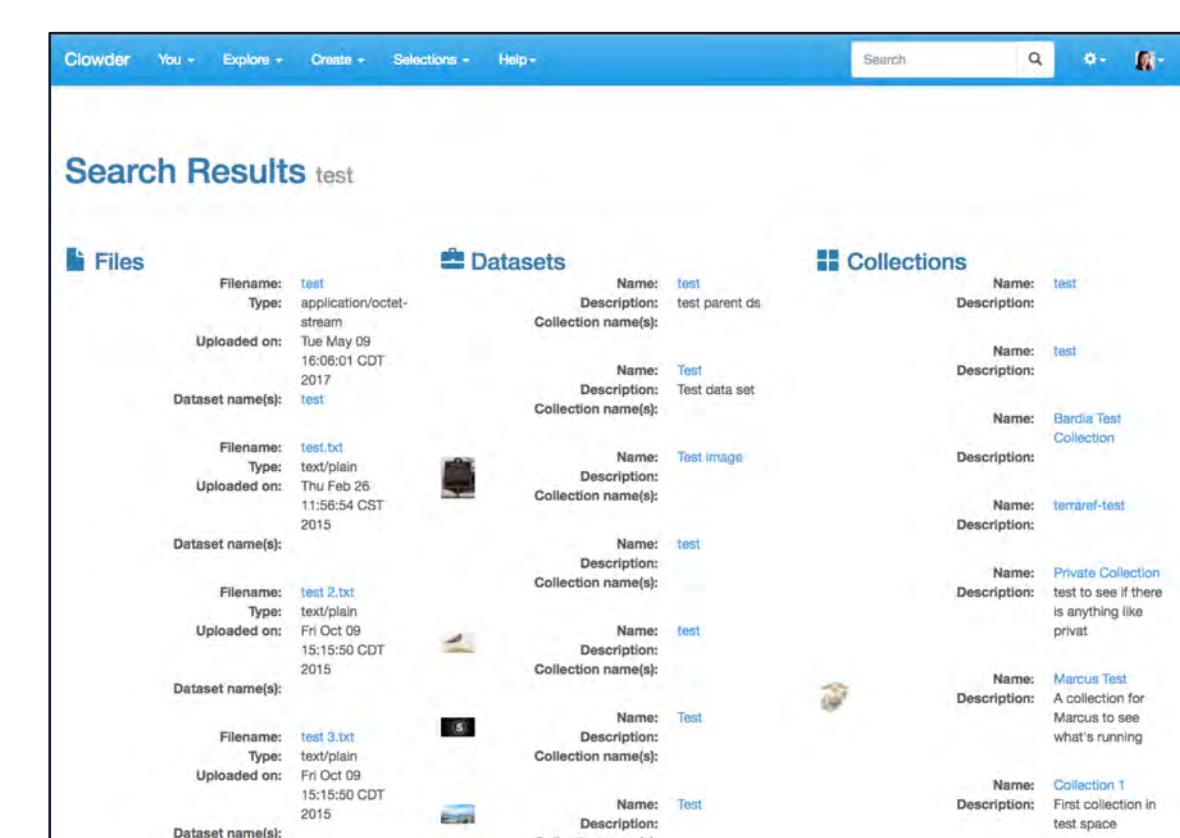
ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

NCSA