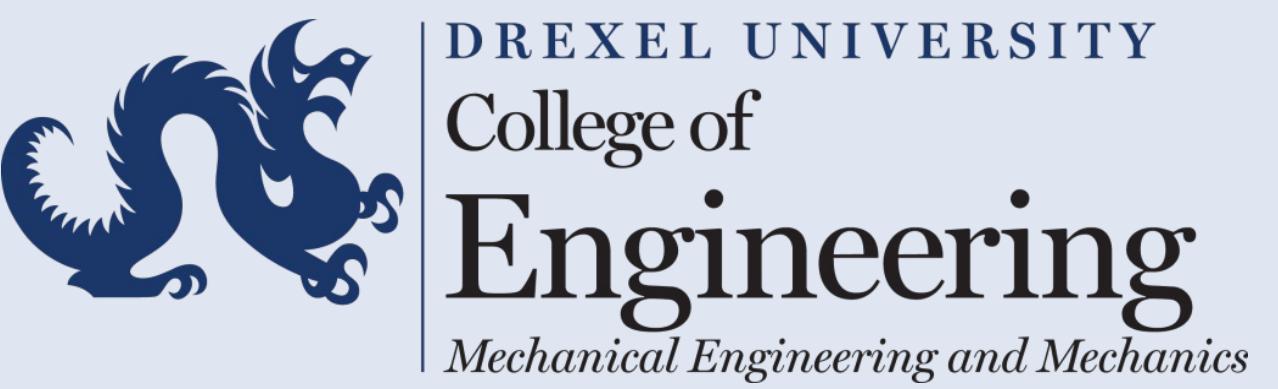


Cyberinfrastructure for Scientific Data Preservation and Image Similarity Search

Yichen Guo^{1,2}, Yifan Zhang³, Julian Goddy², Kio Polson⁴, Kaushik Jagini³, Joshua Brown⁵, Marina Potapova⁶, Chad Peiper⁴, Jane Greenberg⁴, Joshua Agar², Jeff Heflin³

¹Lehigh University Department of Materials Science and Engineering; ²Drexel University Department of Mechanical Engineering and Mechanics; ³Lehigh University Department of Computer Science and Engineering; ⁴Drexel University College of Computing and Informatics; ⁵Oak Ridge National Laboratory Data Lifecycles and Scalable Workflows Group; Department of Biodiversity, ⁶Earth and Environmental Science Academy of Natural Sciences, Drexel University



Most Data is Underanalyzed



- Data analysis takes much longer than acquisition → Analysis takes weeks-months
- Data is generally only accessible by originator

Data is not FAIR



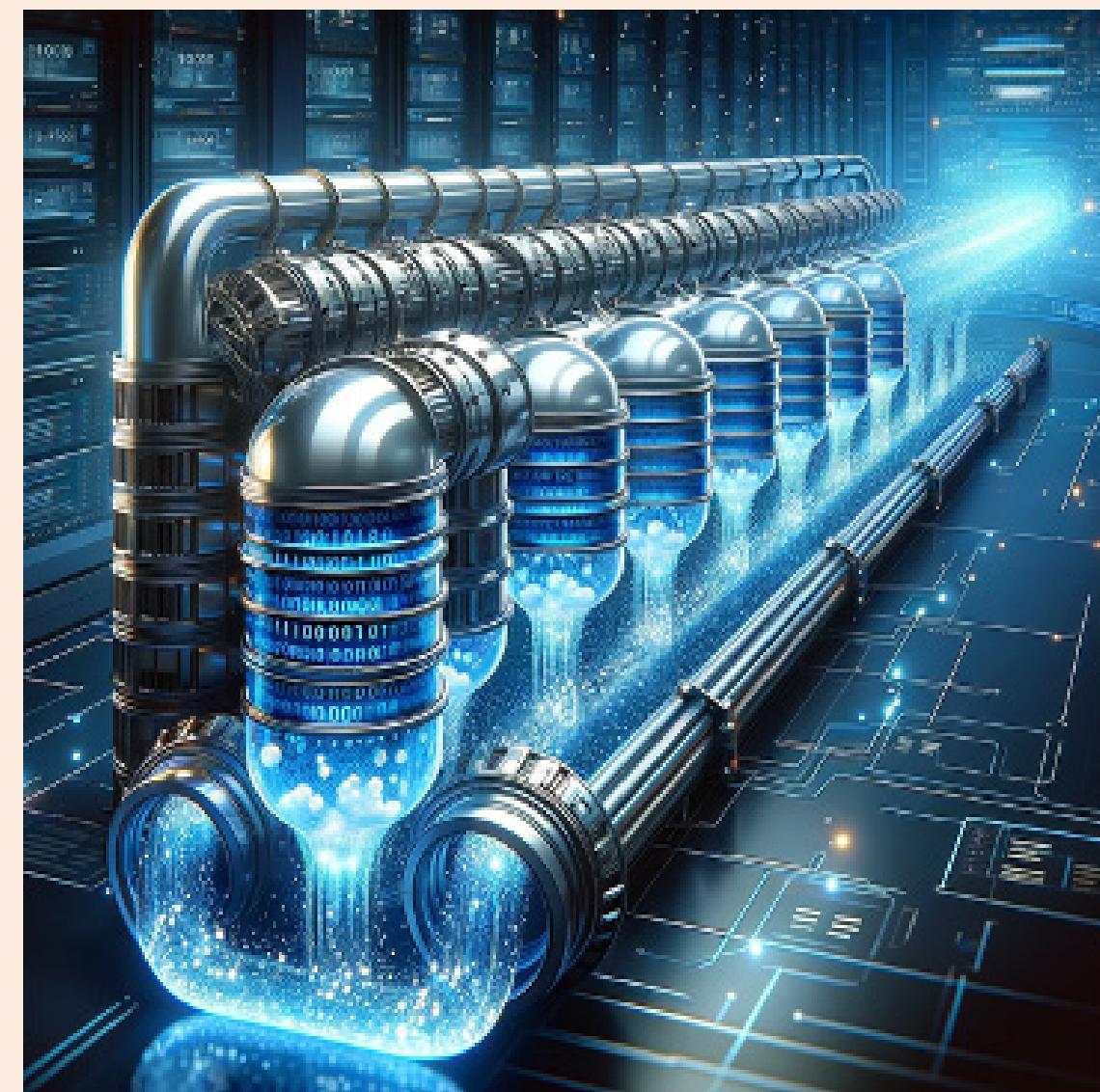
- Science is distributed; it is rare that data is collated → Most data is saved in folders in local file systems
- Sharing between institutions is challenging

Cyberinfrastructure Challenges for Experimental Sciences



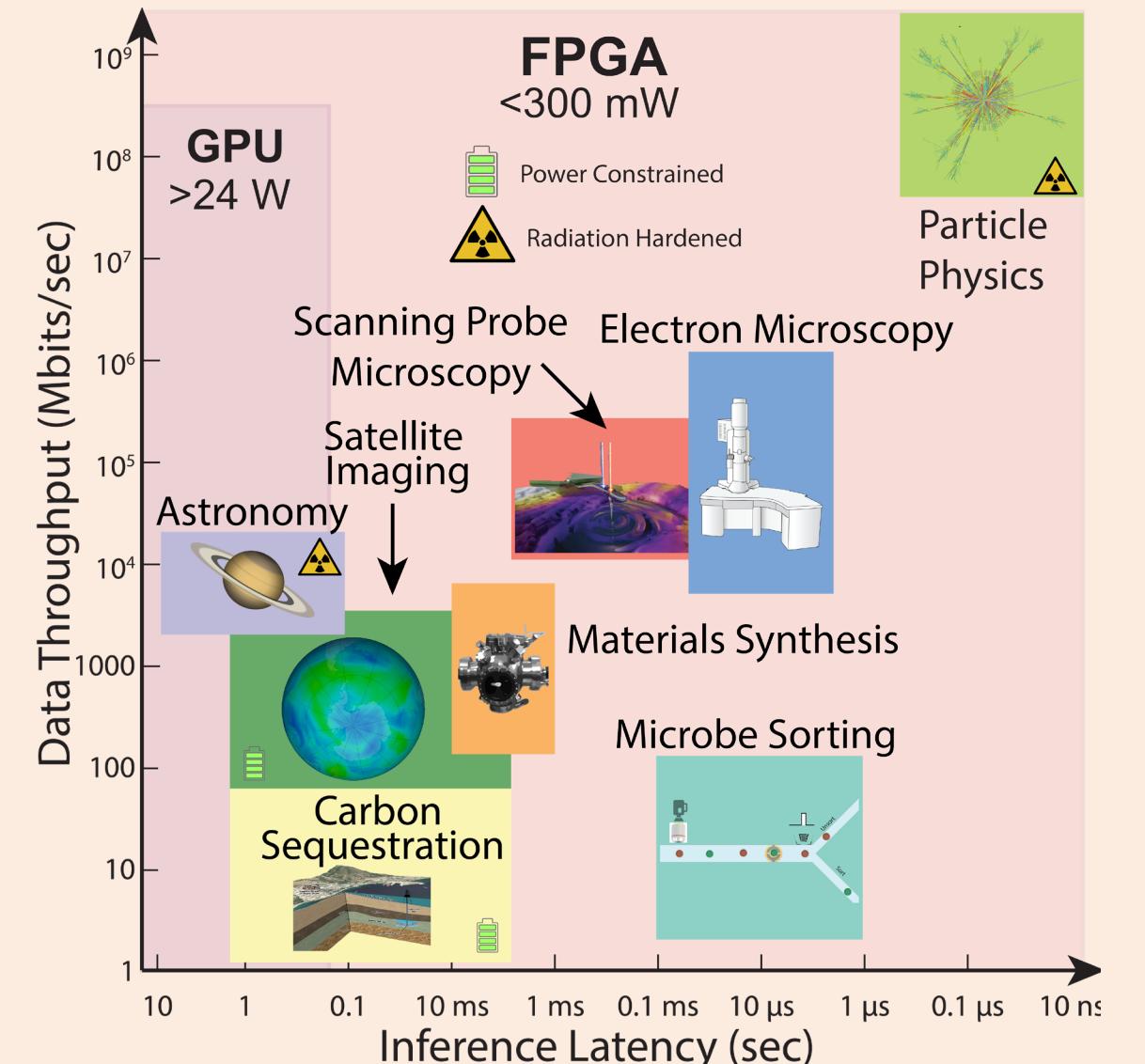
- Experimental scientists have training for functional computational literacy → Minimal support for software development
- Software contributions are undervalued

Computation is Rarely Highly-Available



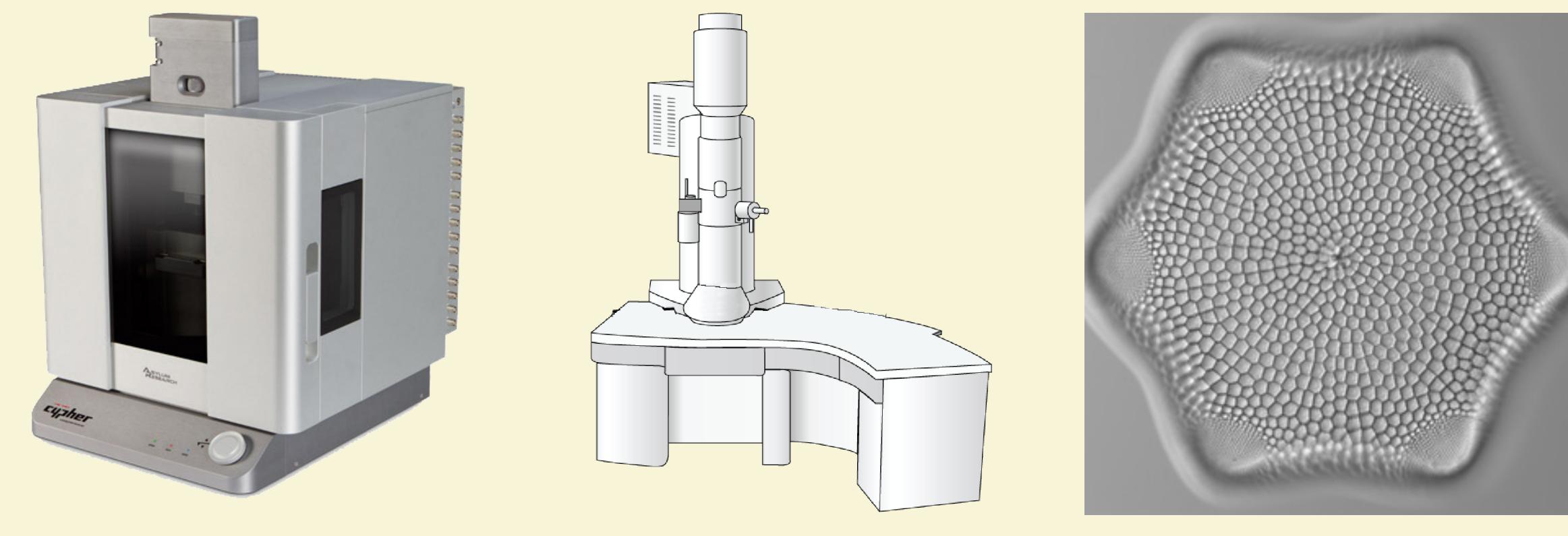
- Compute infrastructure is designed for simulations not experiments → Experiments cannot wait in a queue
- Need for high-availability infrastructure

Non-Deterministic Computational Latency



- Experimentalists rarely deploy deterministic low-latency computation → excluding dynamic process control
- Software, algorithm, hardware codesign

Scientific Data Ingestion



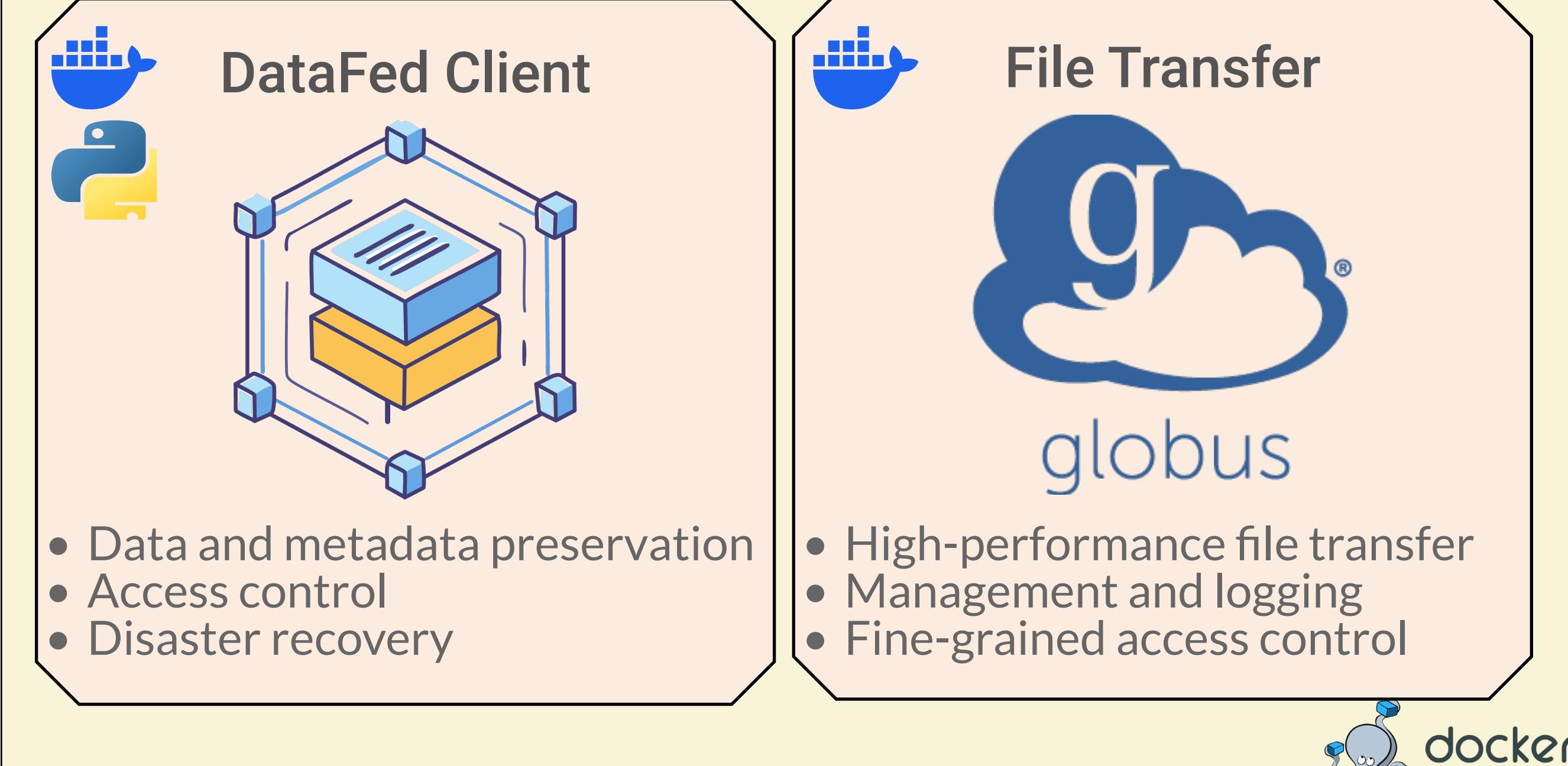
Atomic Force Microscopy Electron Microscopy Diatom Herbarium

Data
↓
Metadata

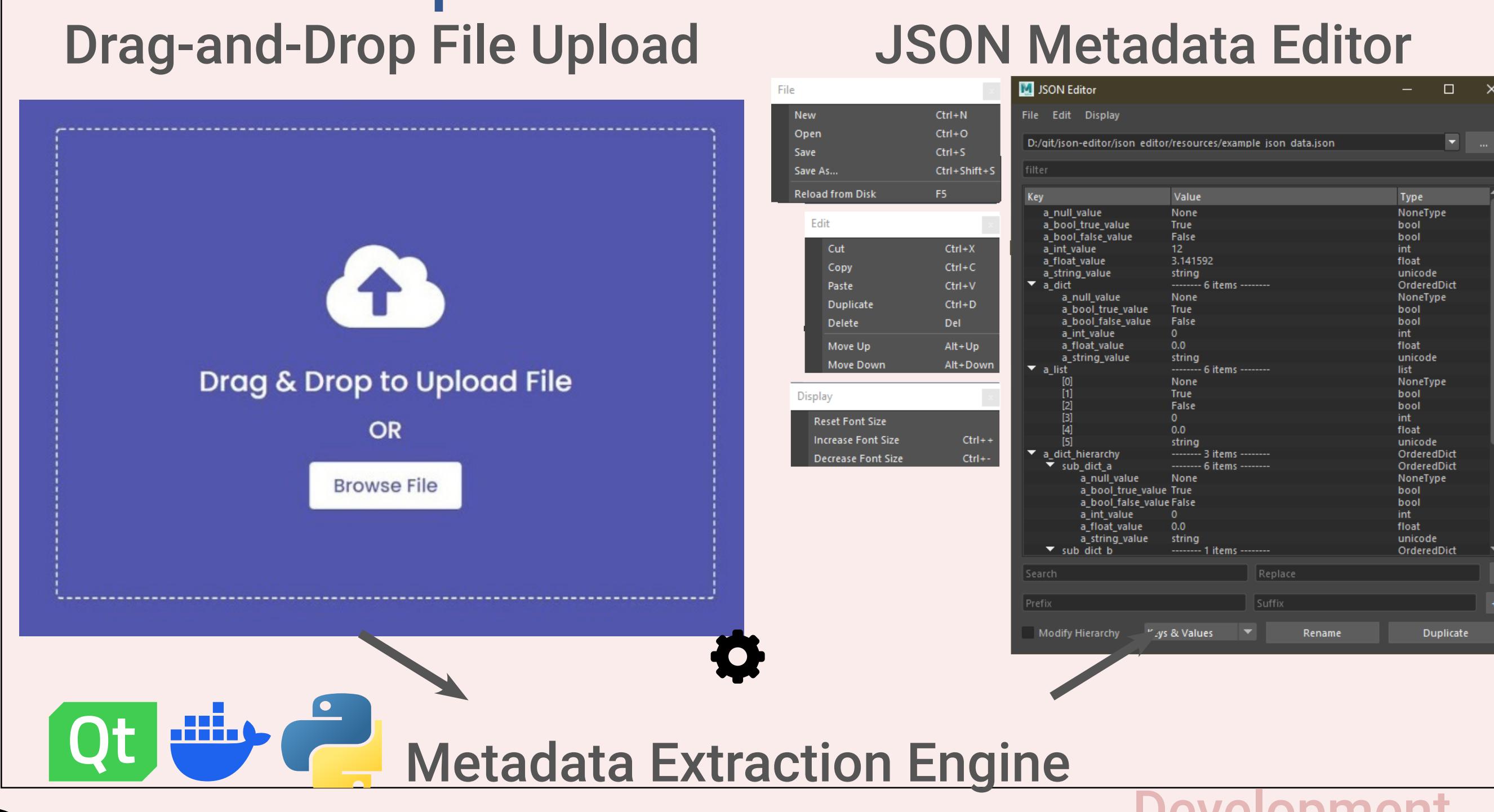
Data
↓
Metadata

Data
↓
Metadata

Composable Containerized Backend



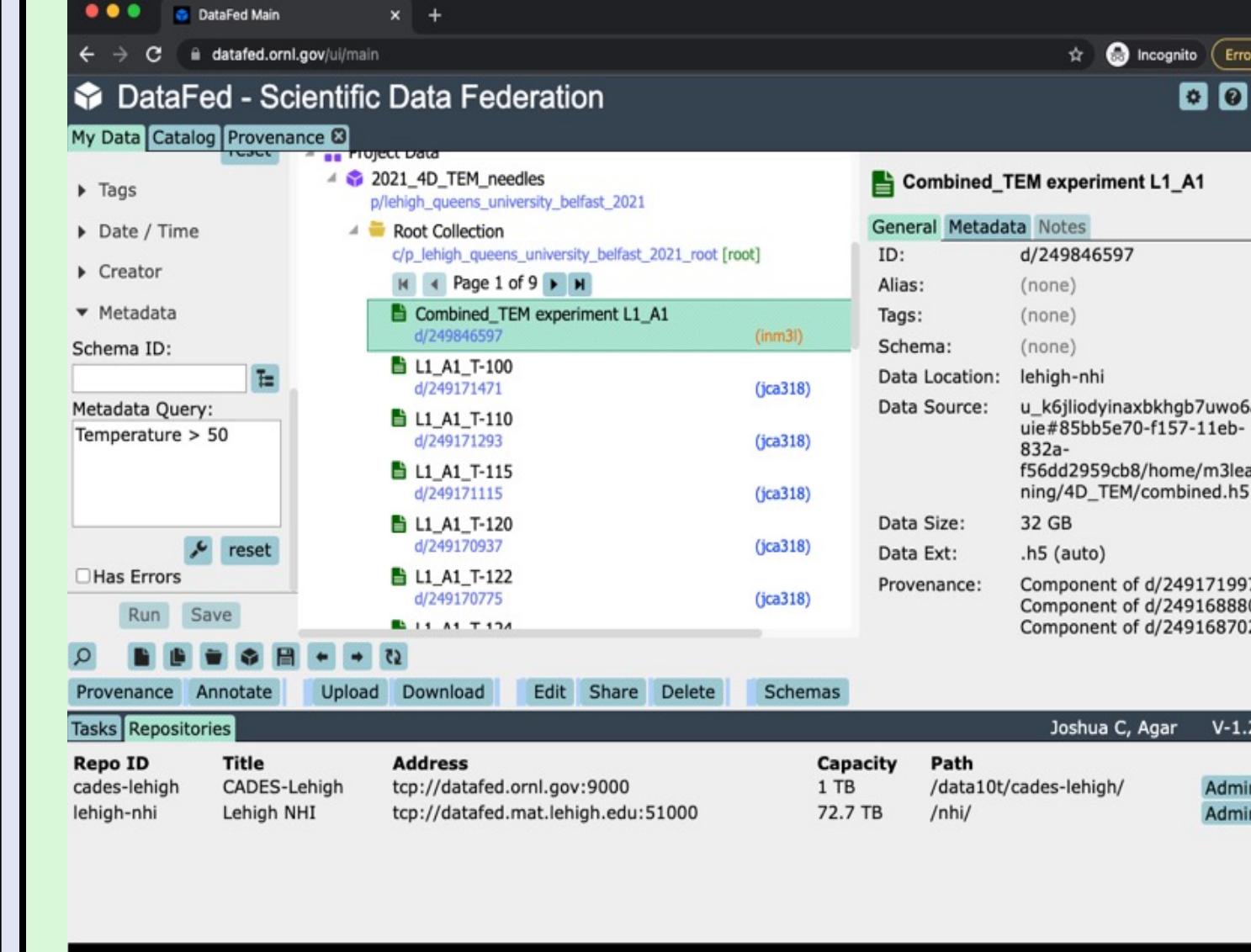
Graphical User Interface



Development

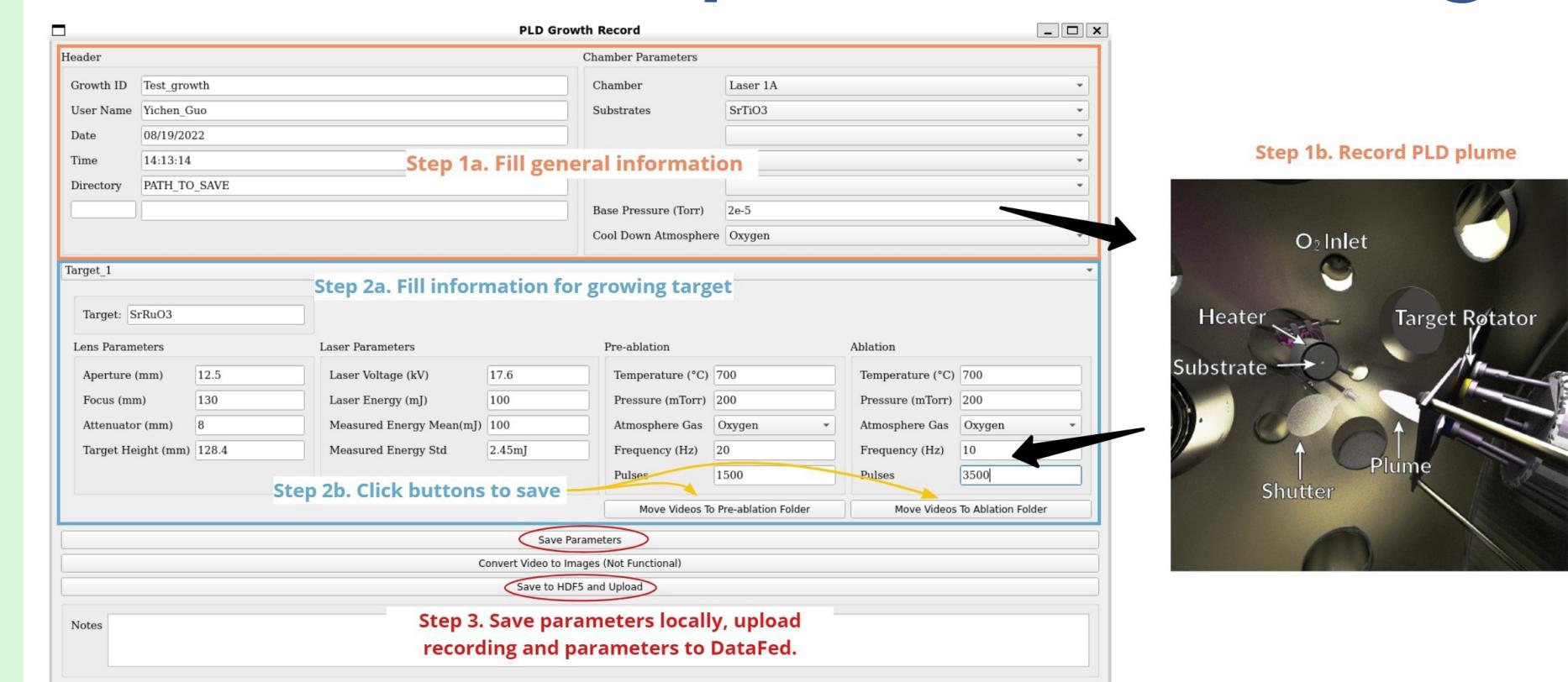
DataFed

Web Interface for Administration and Search



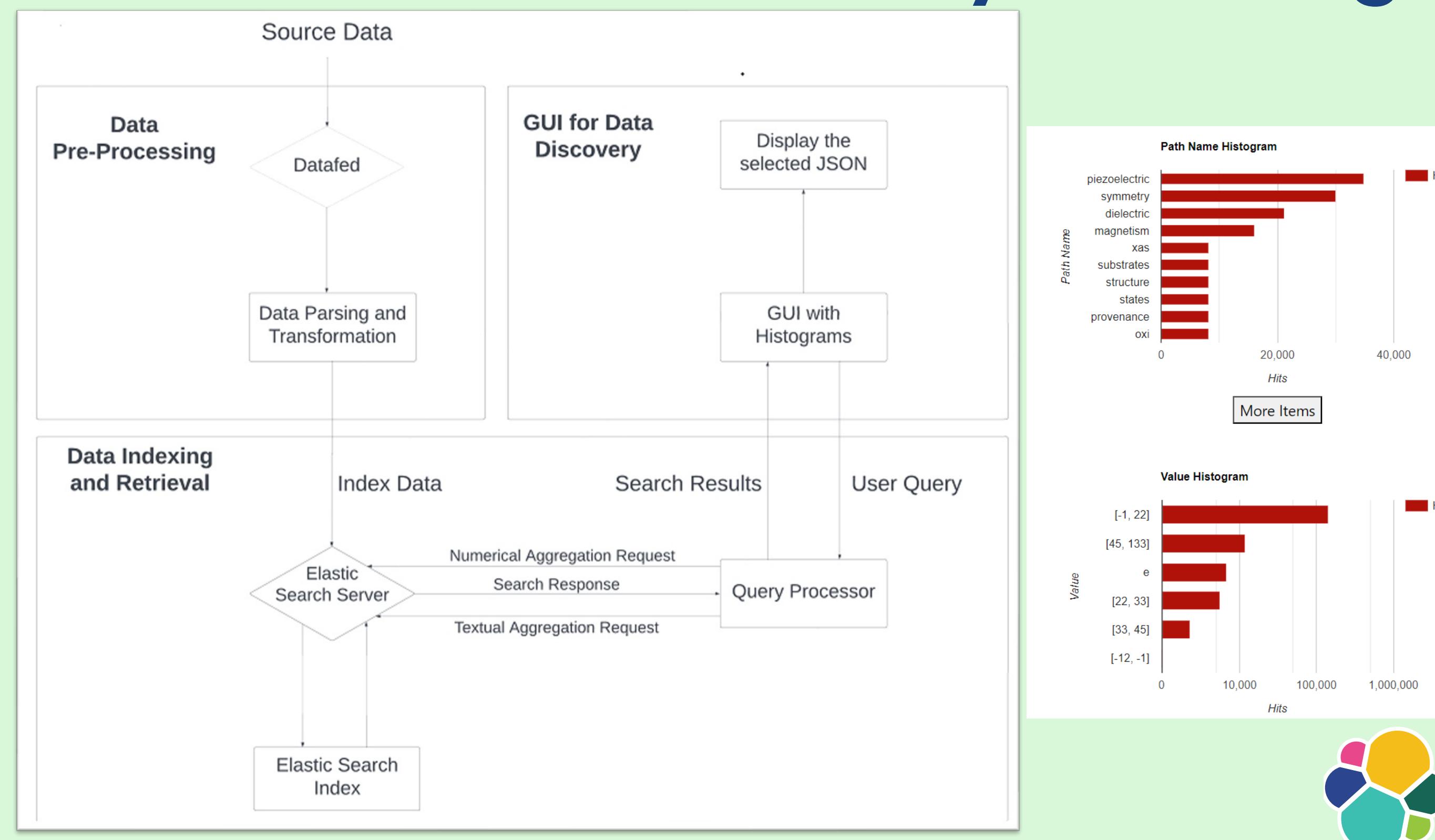
- Federated scientific data management system
- Read, write, and admin control at the user and group level
- Automated file collation and transfer via Globus
- Allows for secure access controlled file transfer between institutional firewalls
- Allows for standard schemes as complex graph relational queries
- Fully functional command line interface and Python API

Python API for Experiment Integration



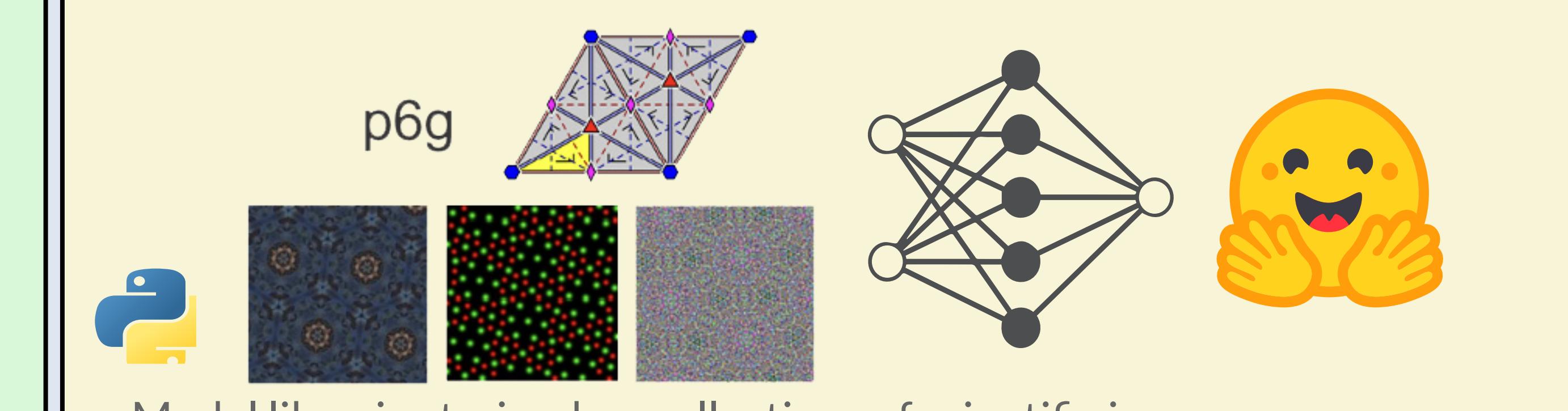
Production

Cell-Centric Discovery Indexing

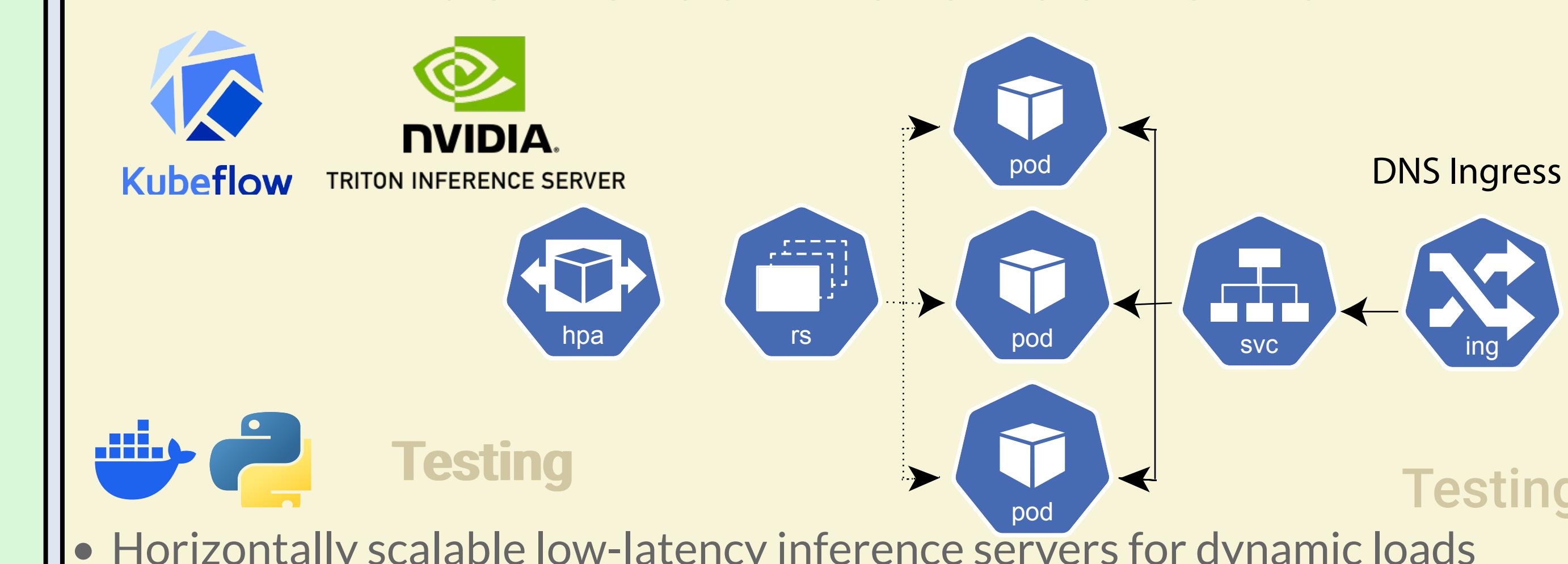


- Novel approach to exploring collections of semi-structured data by extending the cell-centric indexing approach
- Easy data access and retrieval without knowledge of schema or data organization
- User-friendly interface for data exploration and retrieval

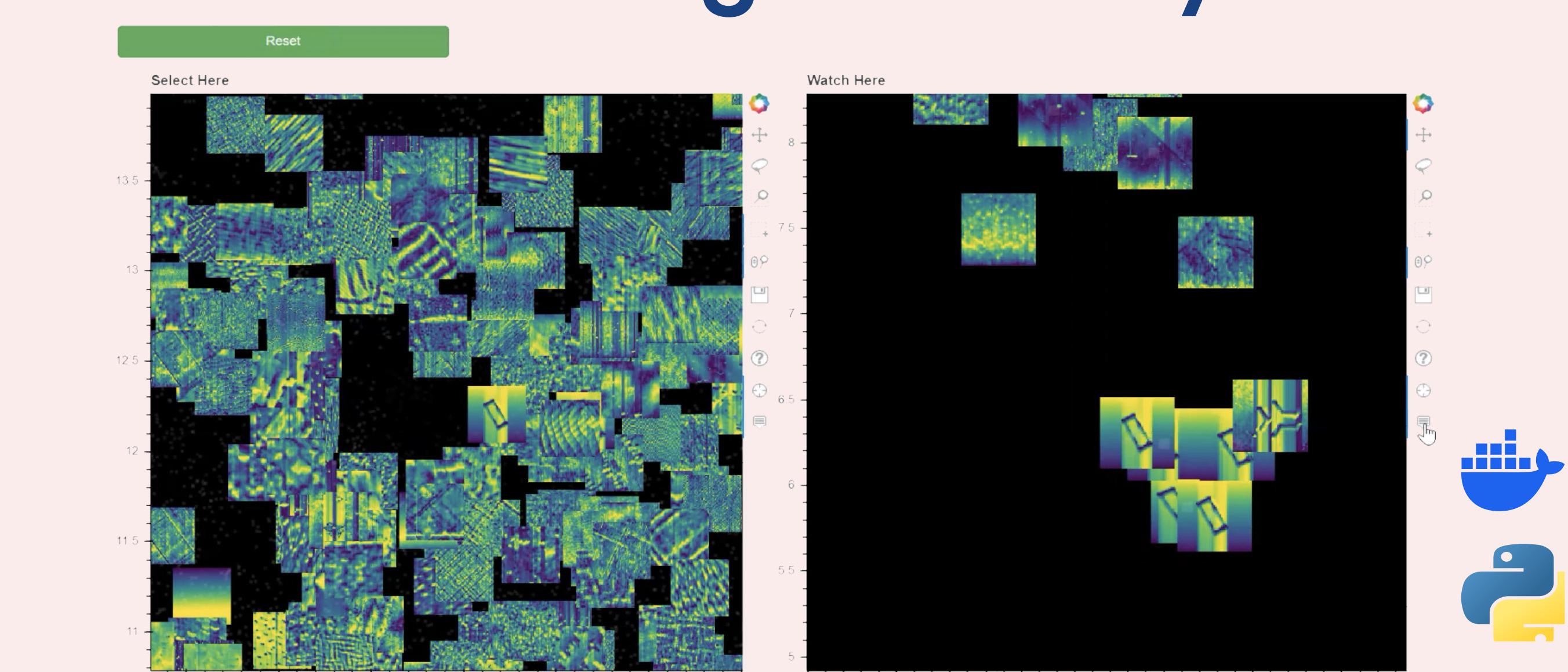
AI Similarity Engine Domain Specific Models



Kubernetes Inference Server



Recursive Image Similarity Search



Funding

Secondary

NSF: MRI: Development of Heterogeneous Edge Computing Platform for Real-Time Scientific Machine Learning (2215789)
NSF: MRI: Development of a Platform for Accessible Data-Intensive Science and Engineering (2320600)
DOE: Real-time Data Reduction Codesign at the Extreme Edge for Science (CHARM) (W911NF-19-2-0119)
Pls: Joshua Agar, Jeff Heflin; 2209135

Primary

Elements: CRISPS; Cell-Centric Recursive Image Similarity Projection Searching