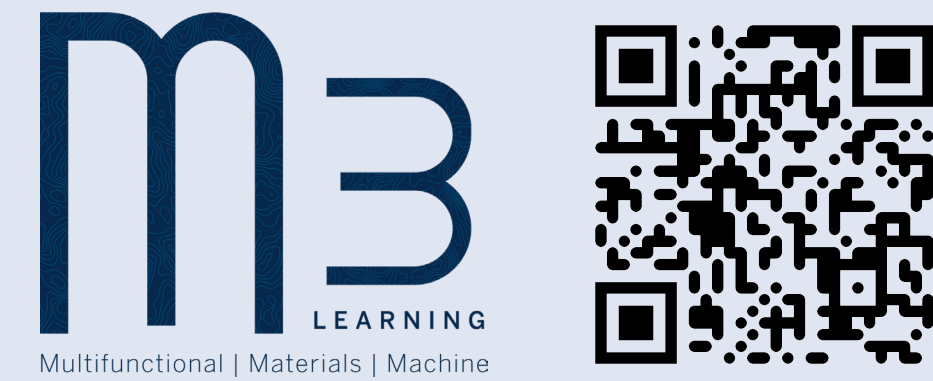
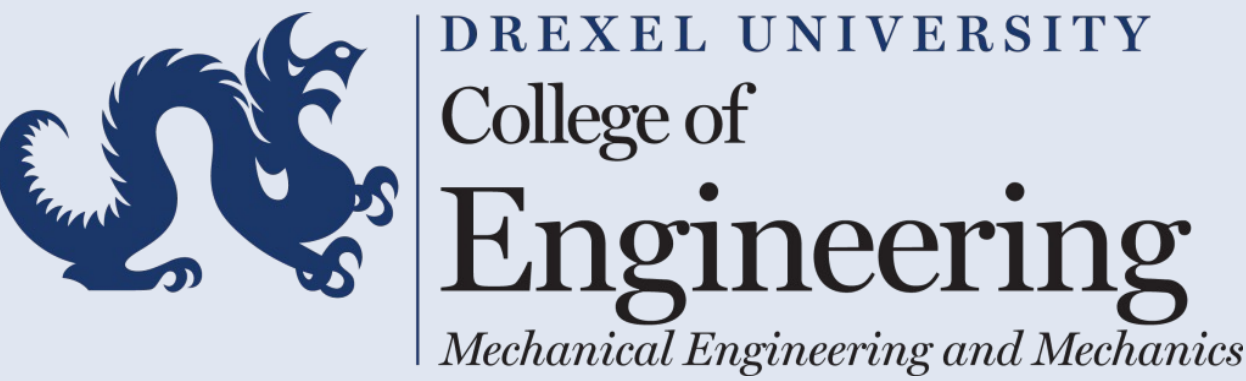


# Cyberinfrastructure for Scientific Data Preservation and Image Similarity Search

Yichen Guo<sup>1,2</sup>, Yifan Zhang<sup>3</sup>, Julian Goddy<sup>2</sup>, Kio Polson<sup>4</sup>, Kaushik Jagini<sup>3</sup>, Joshua Brown<sup>5</sup>, Marina Potapova<sup>6</sup>, Chad Peiper<sup>4</sup>, Jane Greenberg<sup>4</sup>, Joshua Agar<sup>2</sup>, Jeff Heflin<sup>3</sup>

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



## Cyberinfrastructure Challenges for Experimental Sciences

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



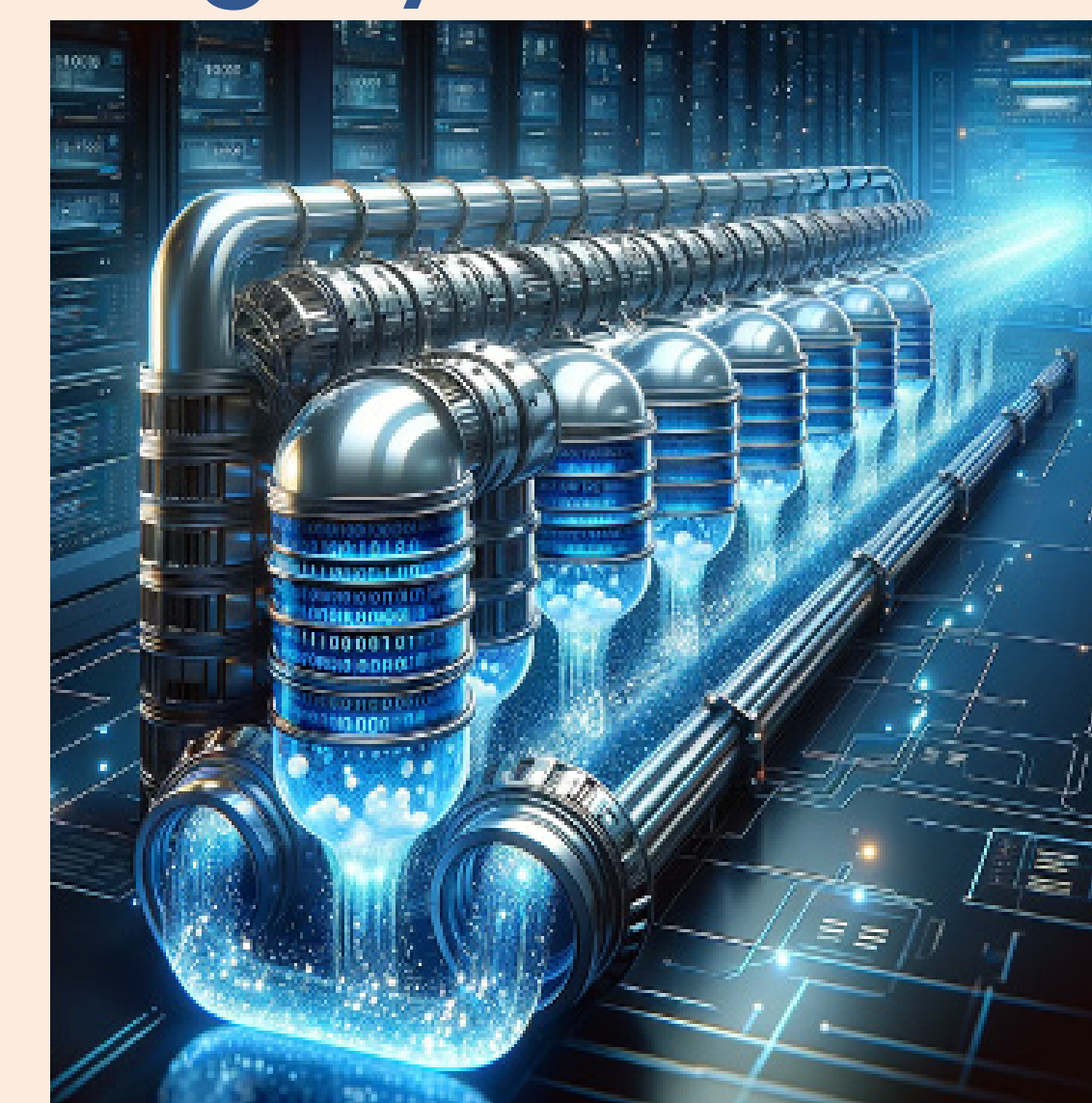
- Science is distributed; it is rare that data is collated → Most data is saved in folders in local file systems
- Sharing between institutions is challenging

Cyberinfrastructure is not Science



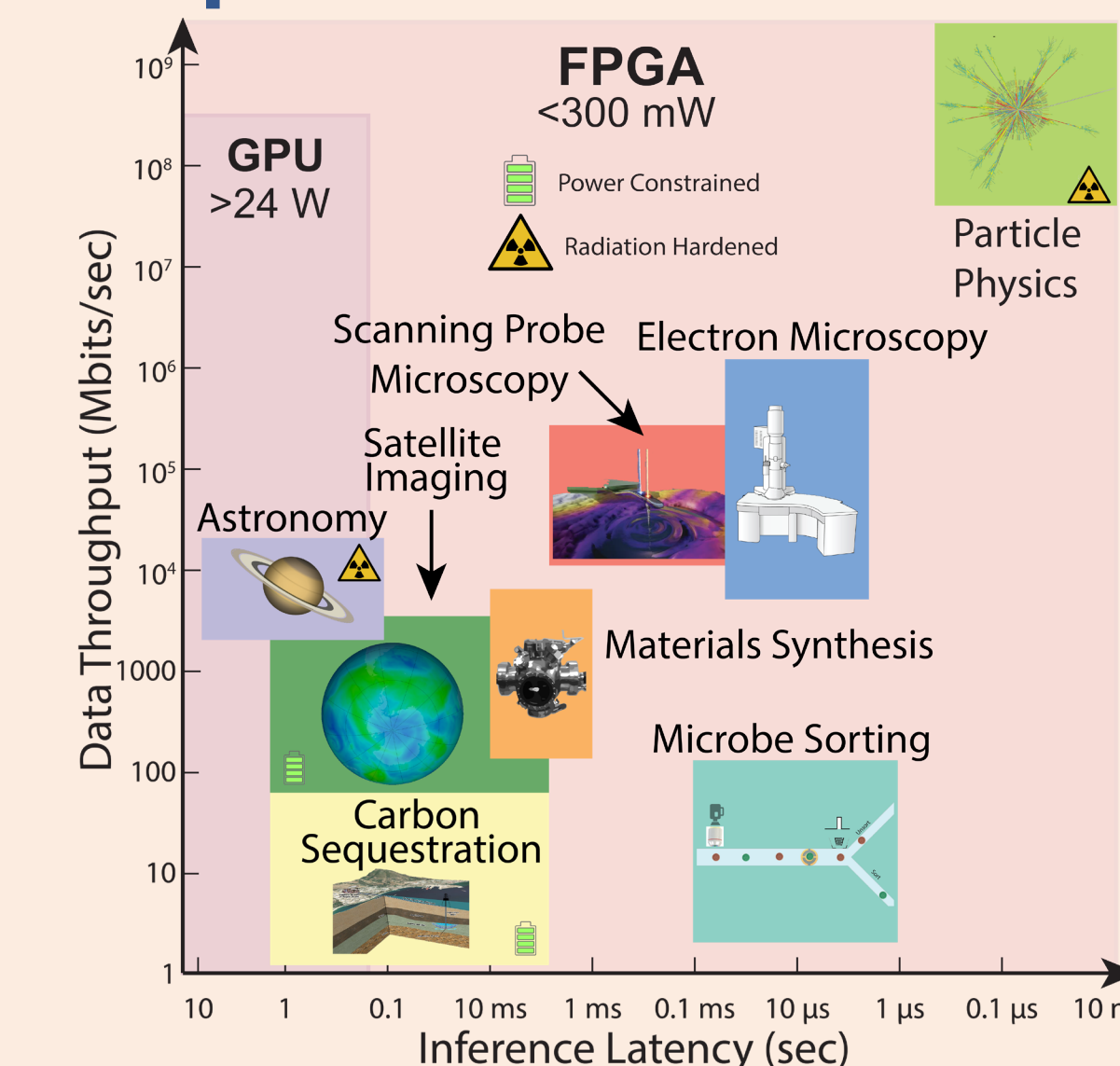
- Experimental scientists have training for functional computational literacy → Minimal support for software development
- Software contributions are undervalued

Computation is Rarely Highly-Available



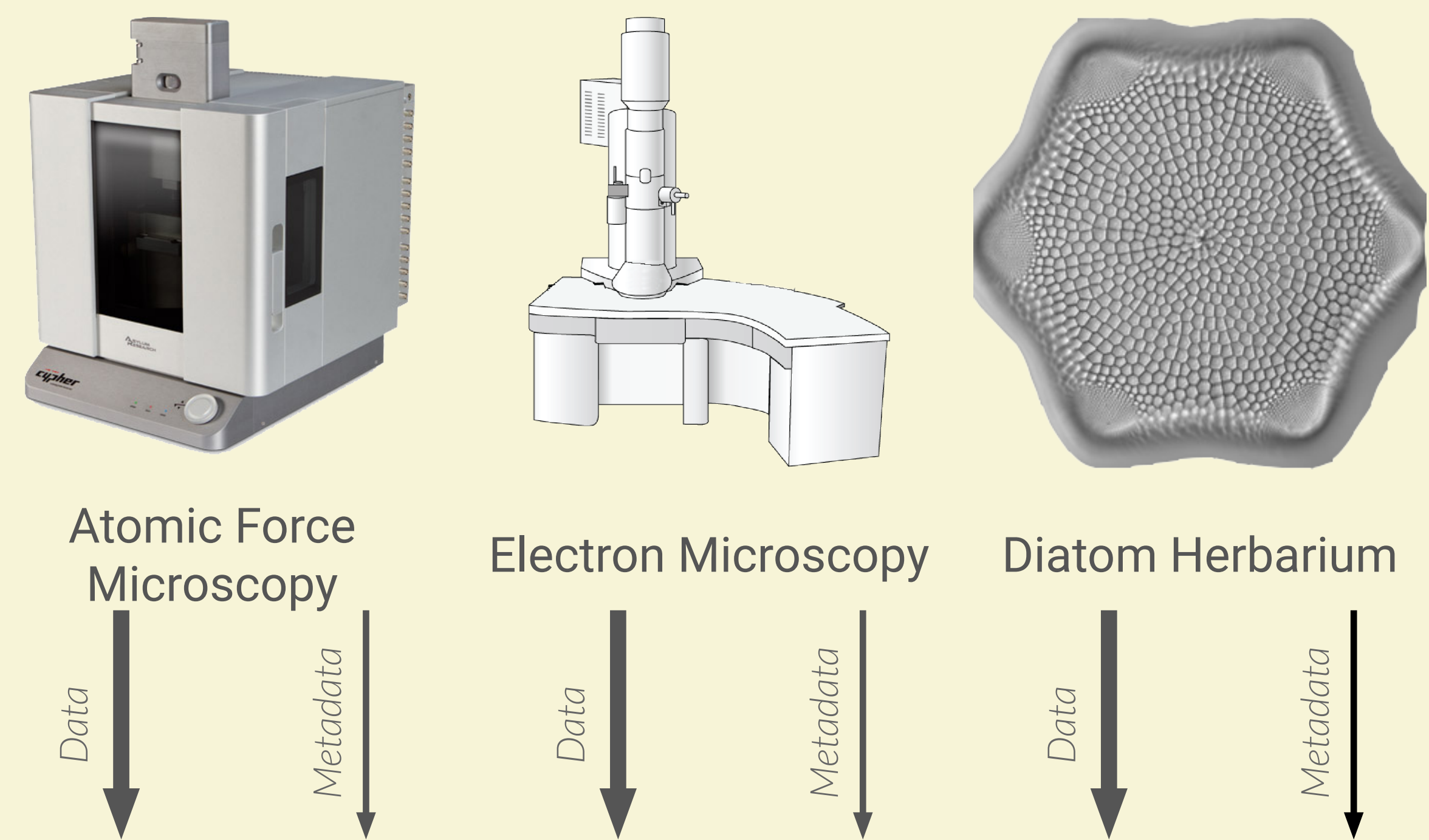
- Compute infrastructure is designed for simulations not experiments → Experiments cannot wait in a queue
- Need for high-availability infrastructure

Non-Deterministic Computational Latency

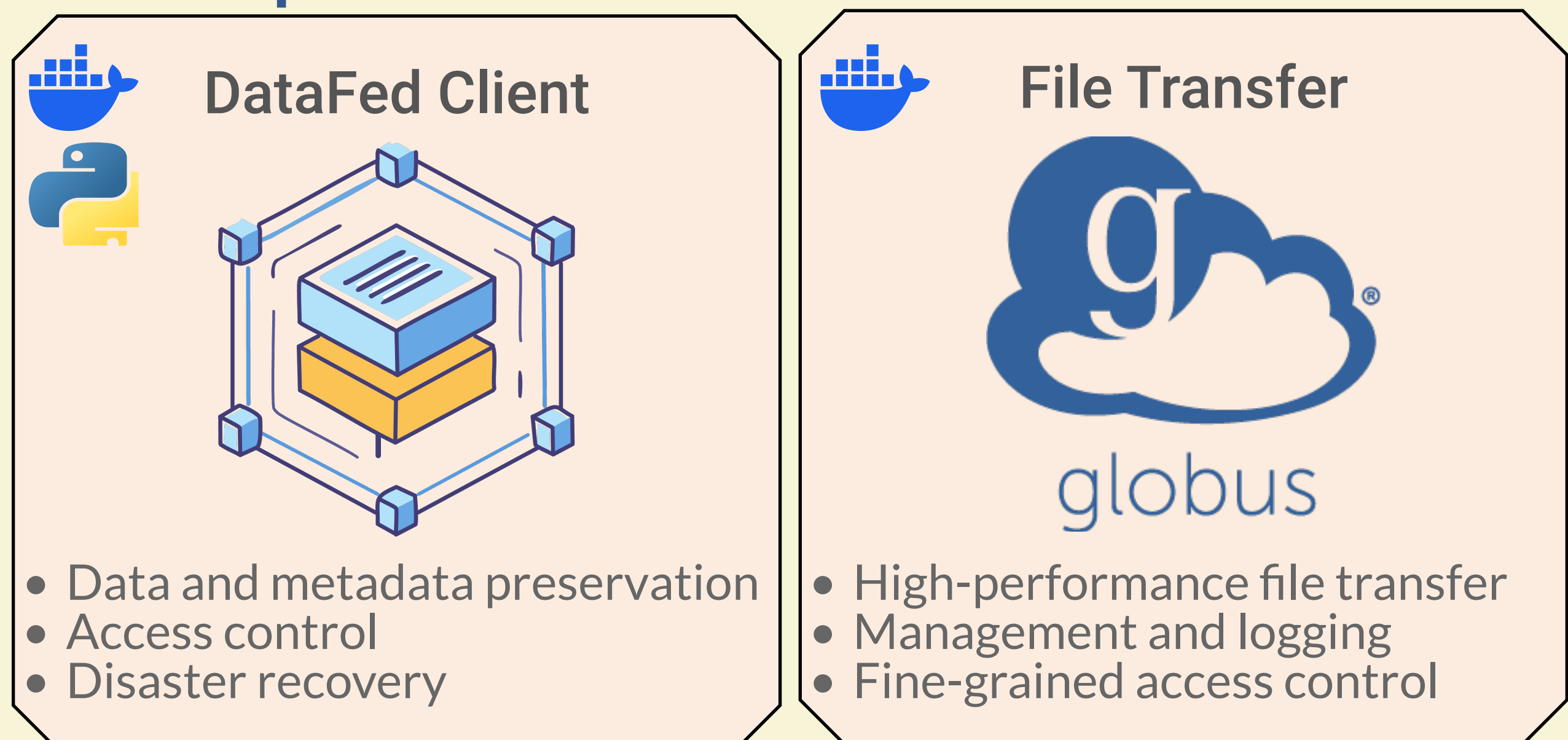


- Experimentalists rarely deploy deterministic low-latency computation → excluding dynamic process control
- Software, algorithm, hardware codesign

### Scientific Data Ingestion

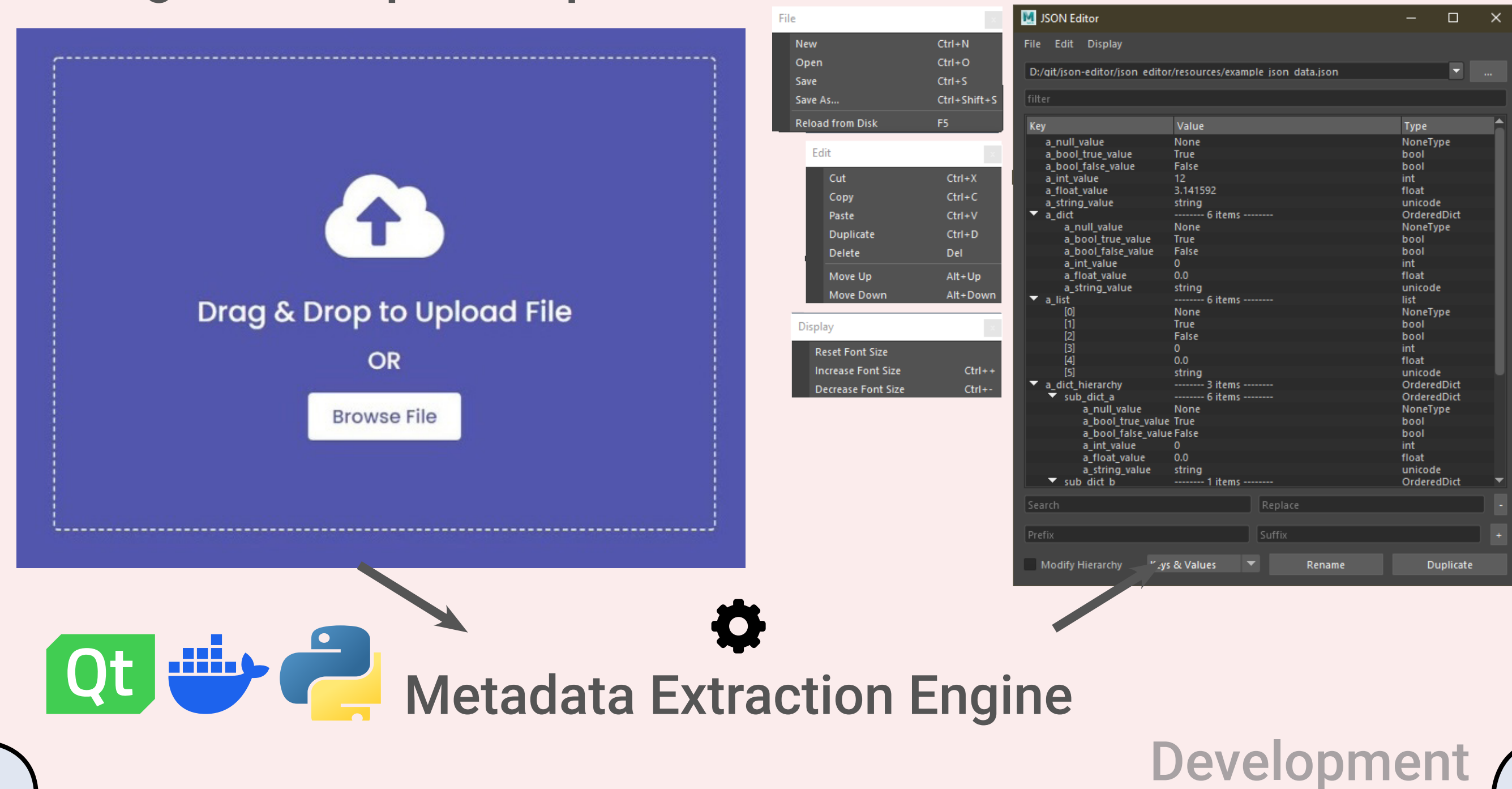


### Composable Containerized Backend



Testing

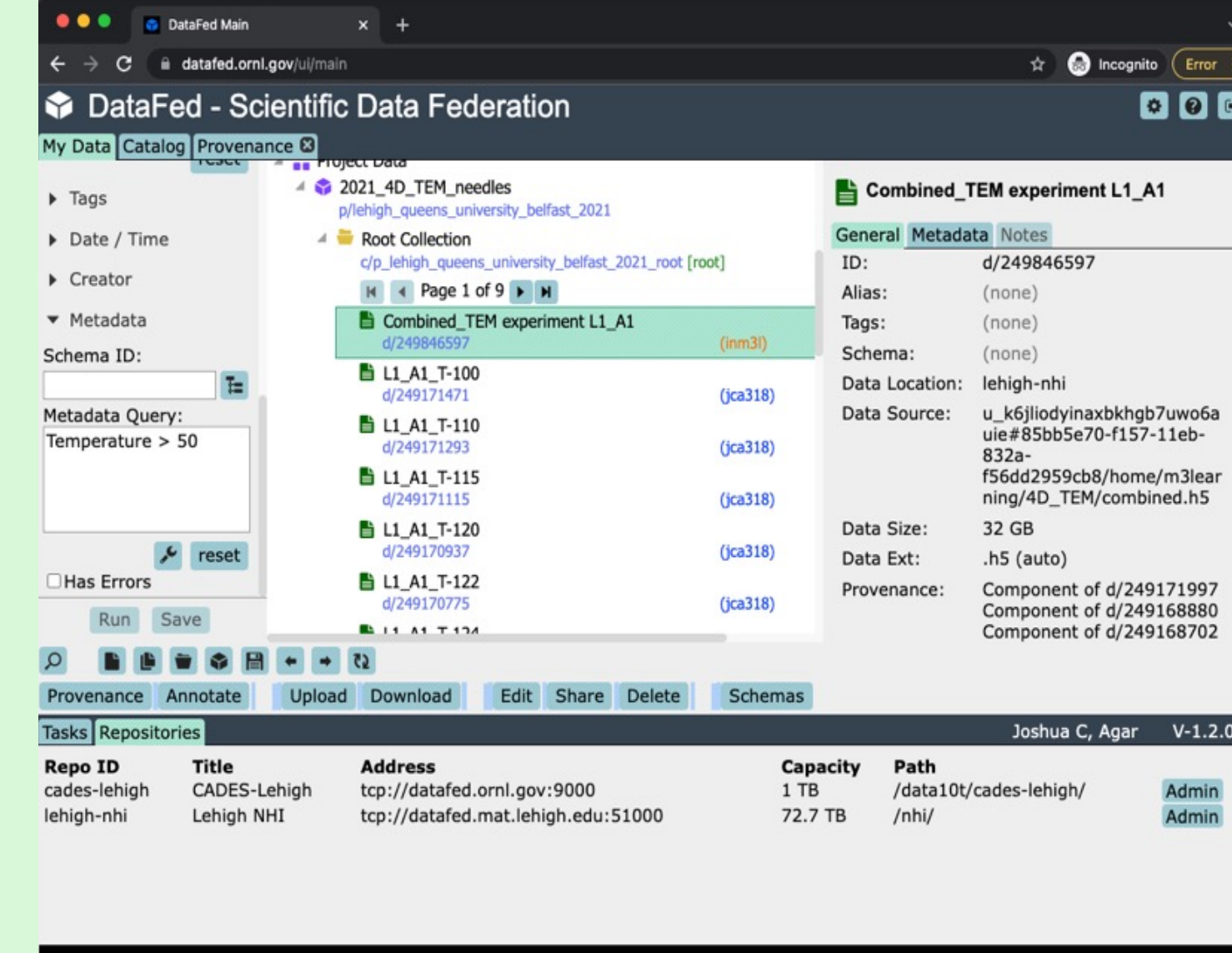
### Graphical User Interface



Development

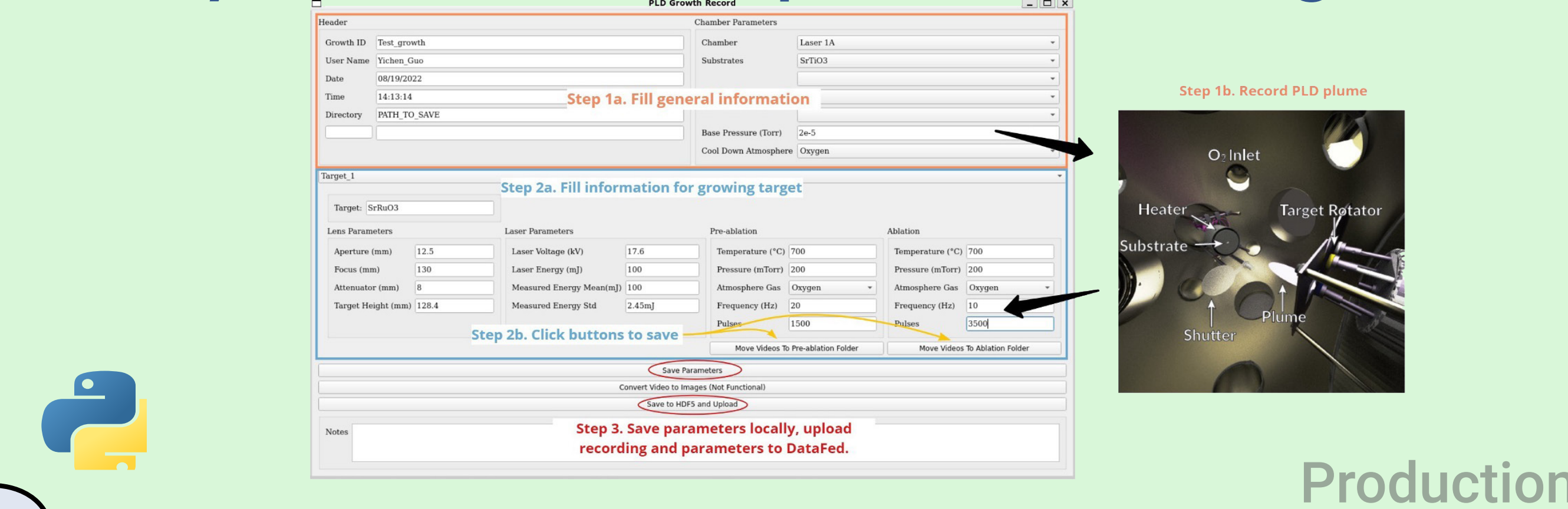
### DataFed

#### 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
- Secure access controlled file transfer between institutional firewalls
- Standard schemes as complex graph relational queries
- Fully functional command line interface and Python API

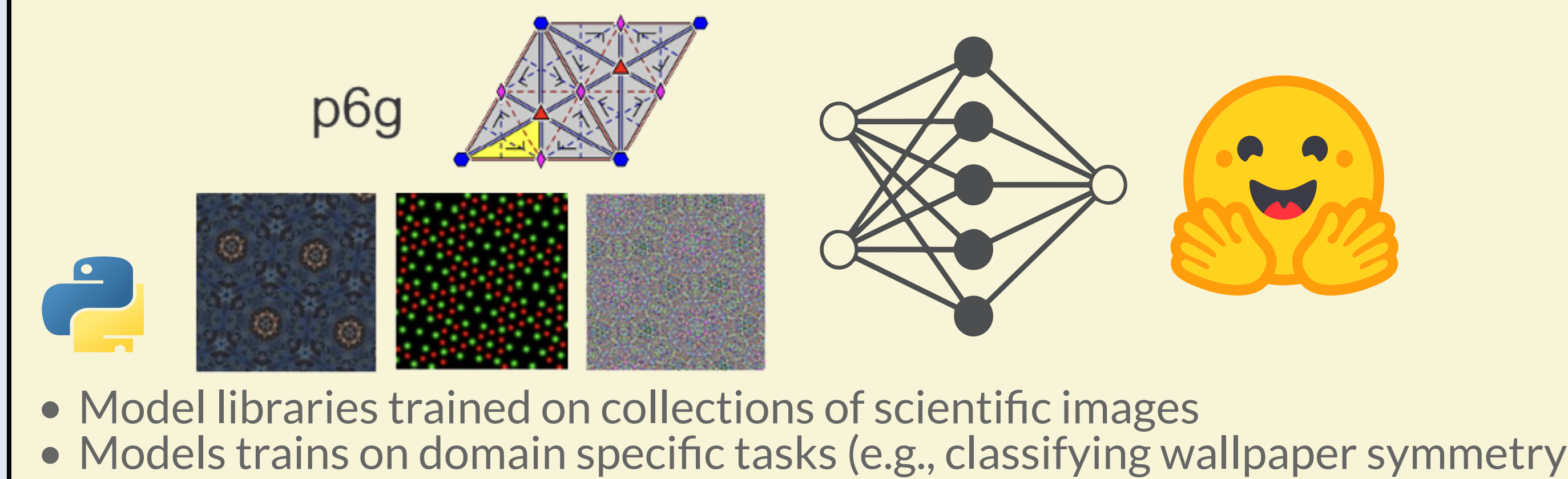
#### Python API for Experiment Integration



Production

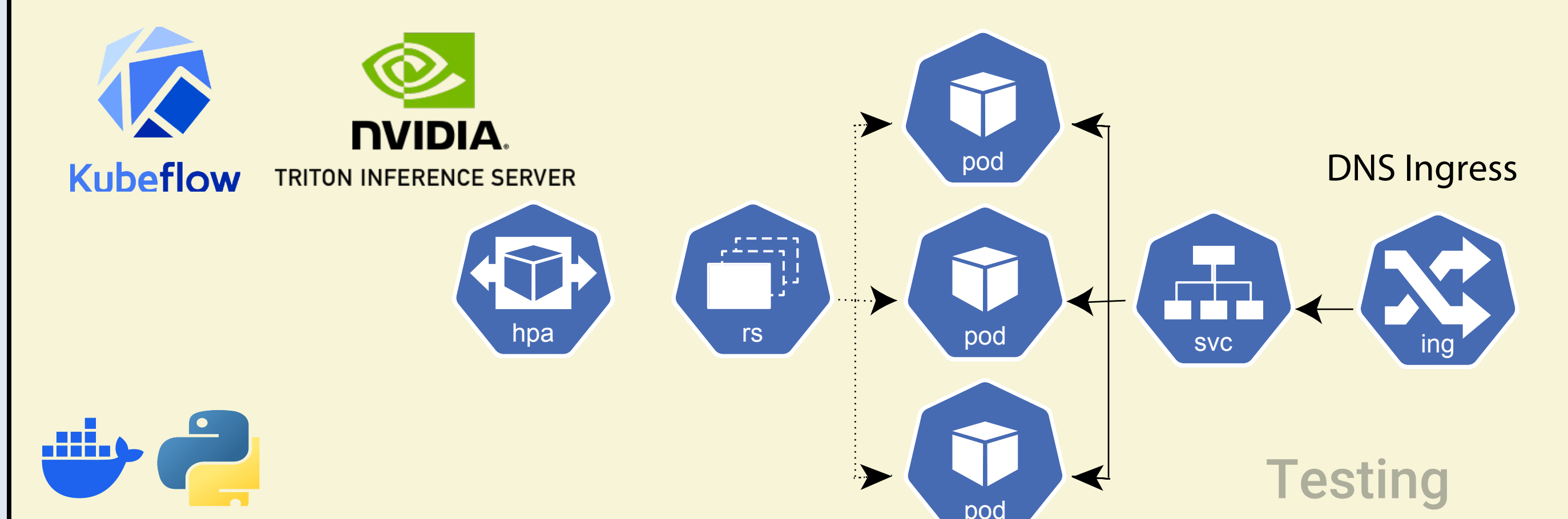
### AI Similarity Engine

#### Domain Specific Models



- Model libraries trained on collections of scientific images
- Models trains on domain specific tasks (e.g., classifying wallpaper symmetry)

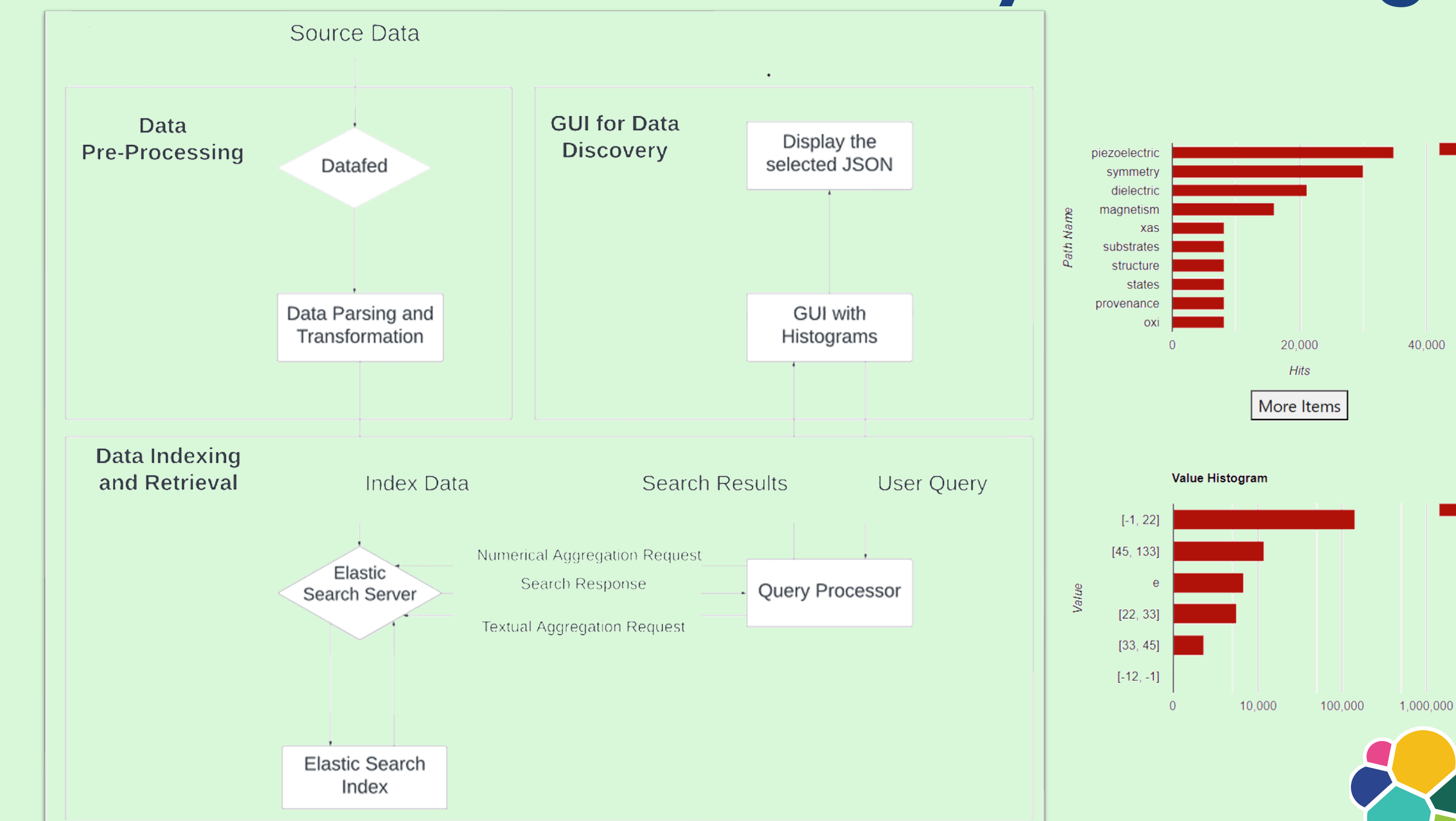
#### Kubernetes Inference Server



- Horizontally scalable low-latency inference servers for dynamic loads
- Allows for rolling model deployment and updates with 0 down time

Testing

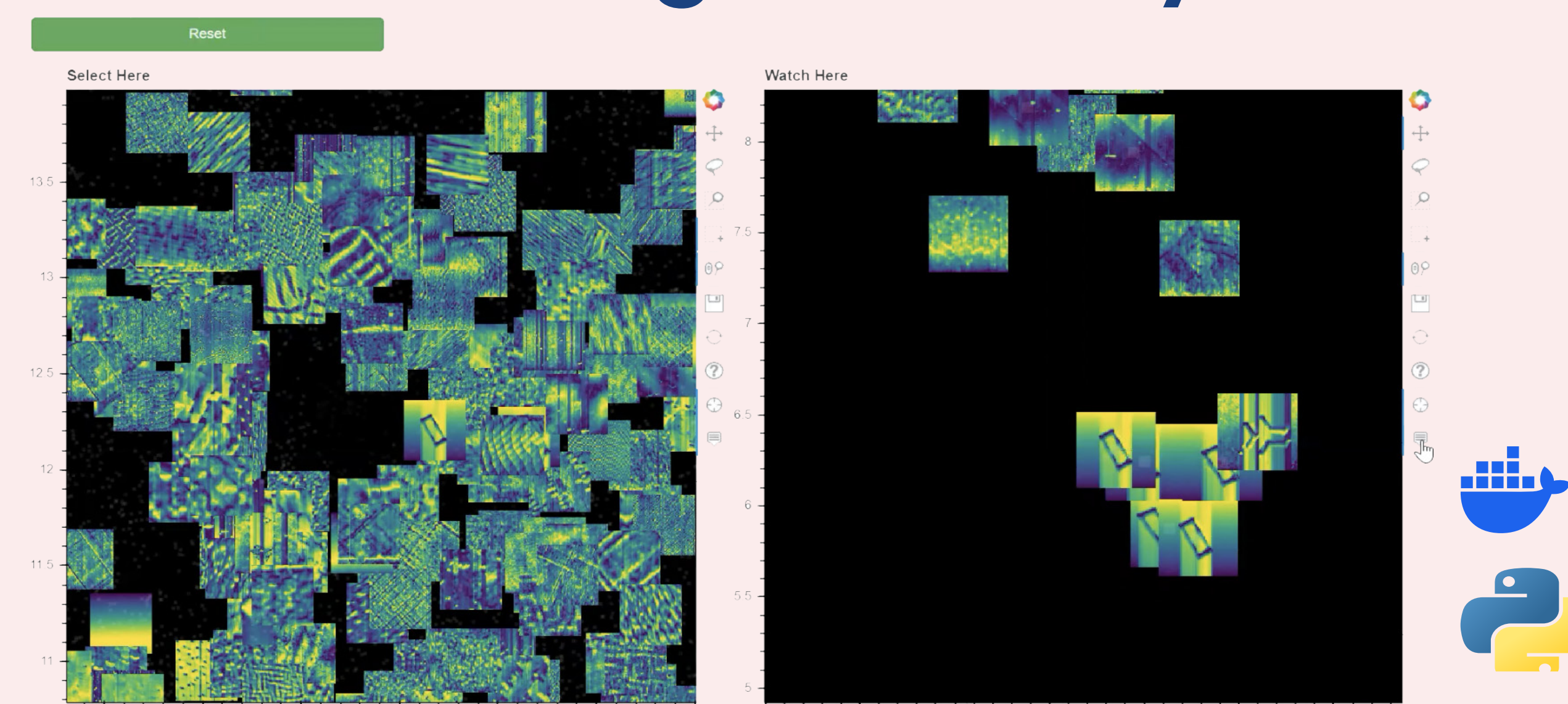
### Cell-Centric Discovery Indexing



- 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

Production

### Recursive Image Similarity Search



- Interactive tool for image similarity projections, and interactive, search
- Facilities discovery without prior knowledge from minimally collated image data

Development

Primary

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



PIs: Joshua Agar, Jeff Heflin; 2209135

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  
ARL: Collaborative for Hierarchical Agile and Responsive Materials (CHARM) (W911NF-19-2-0119)